

Microsoft Research

Each year Microsoft Research hosts hundreds of influential speakers from around the world including leading scientists, renowned experts in technology, book authors, and leading academics, and makes videos of these lectures freely available.

2013 © Microsoft Corporation. All rights reserved.

New methods for the analysis of genome variation data

Richard Durbin

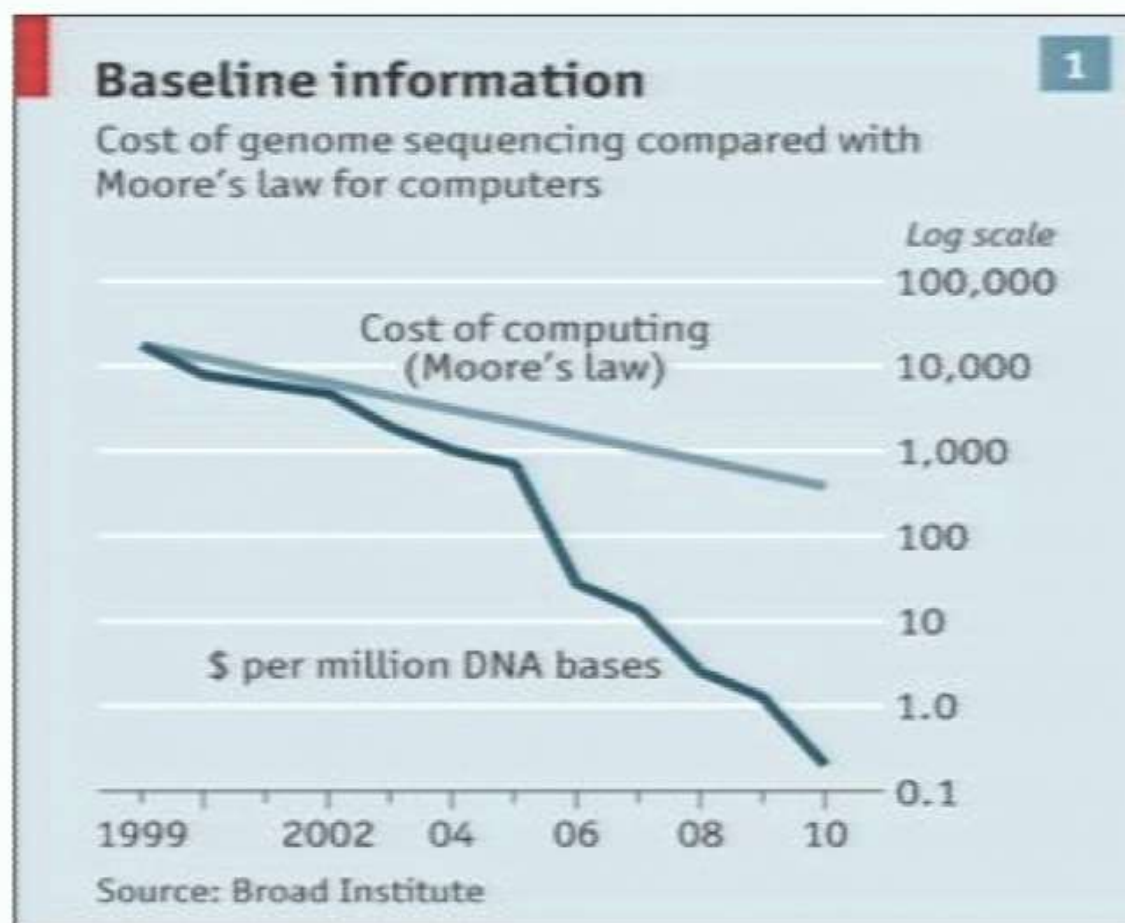
Wellcome Trust Sanger Institute

rd@sanger.ac.uk

The era of sequencing genomes

	Size (Mb)	Genes			Completion date
<i>H. influenzae</i>	2	1,700	1/1kb	Bacterium	1995
Yeast	13	6,000	1/2kb	Eukaryotic cell	1996
Nematode	100	18,000	1/6kb	Animal	1998
Human	3000	20,000	1/150kb	Mammal	2000/3

Between 2000 and 2010 DNA sequencing costs dropped by five orders of magnitude



This scaling, now generating PB data per year, challenges data analysis

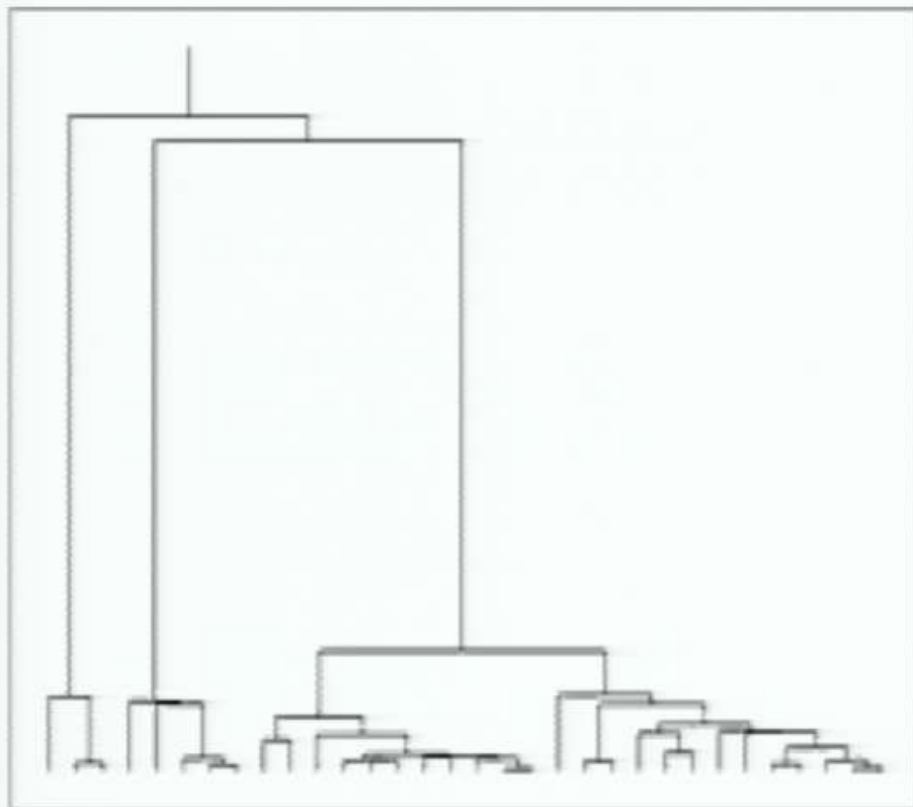
We can now study genetic variation directly by sequencing

- 2000: First human genome ~\$1 billion
- 2009: ~200 human genomes sequenced
 - Human genome \$50k
- 2012/13: Thousands of genomes for research
 - Human genome <\$5k
- 2015/16: Human genome < \$1k
 - Millions of human genomes for clinical usage, research, personal interest

We need to use the structure in the data for representation and inference

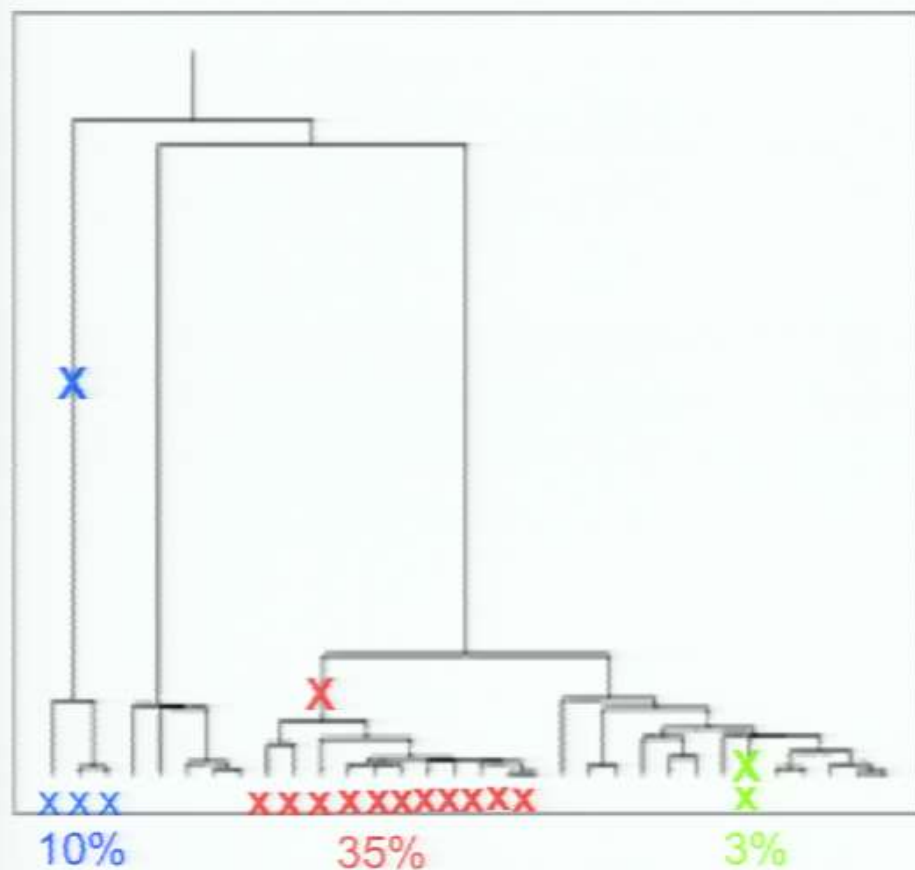
DNA sequences are related by evolution

DNA sequences are related by evolution



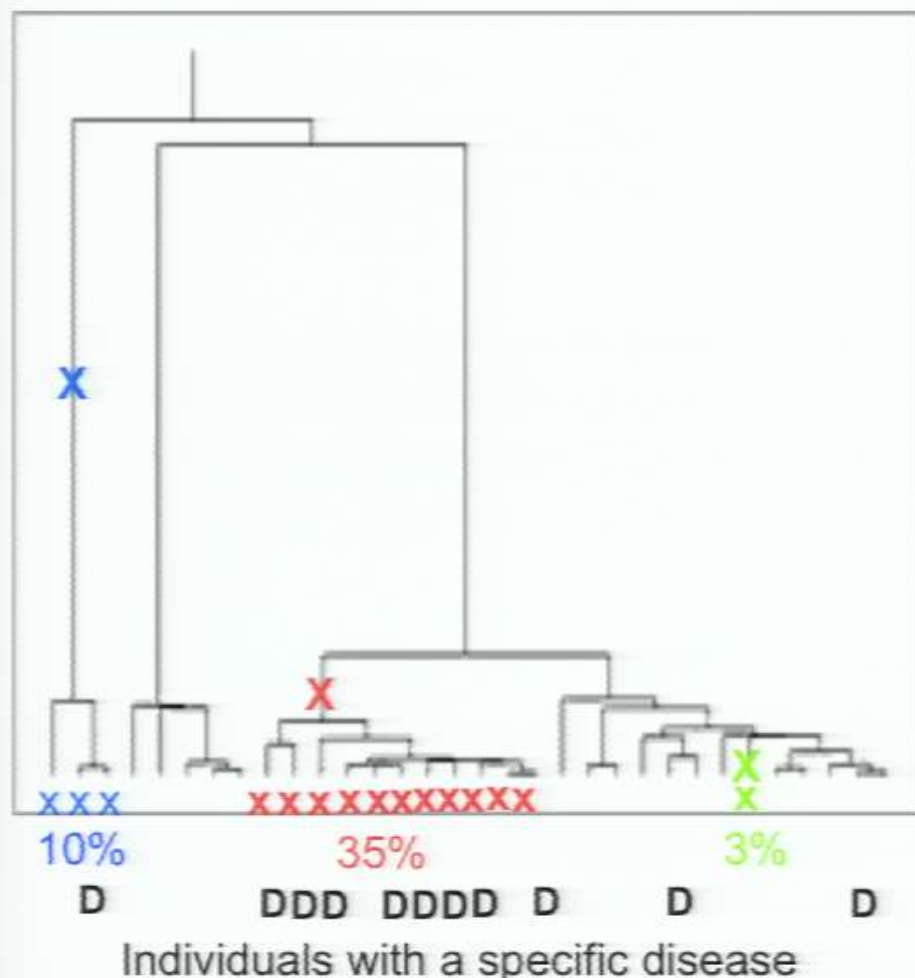
- At each point in the genome we are descended from a common ancestor

DNA sequences are related by evolution



- At each point in the genome we are descended from a common ancestor
- Mutations since the common ancestor give rise to genetic variants shared by the descendants

DNA sequences are related by evolution

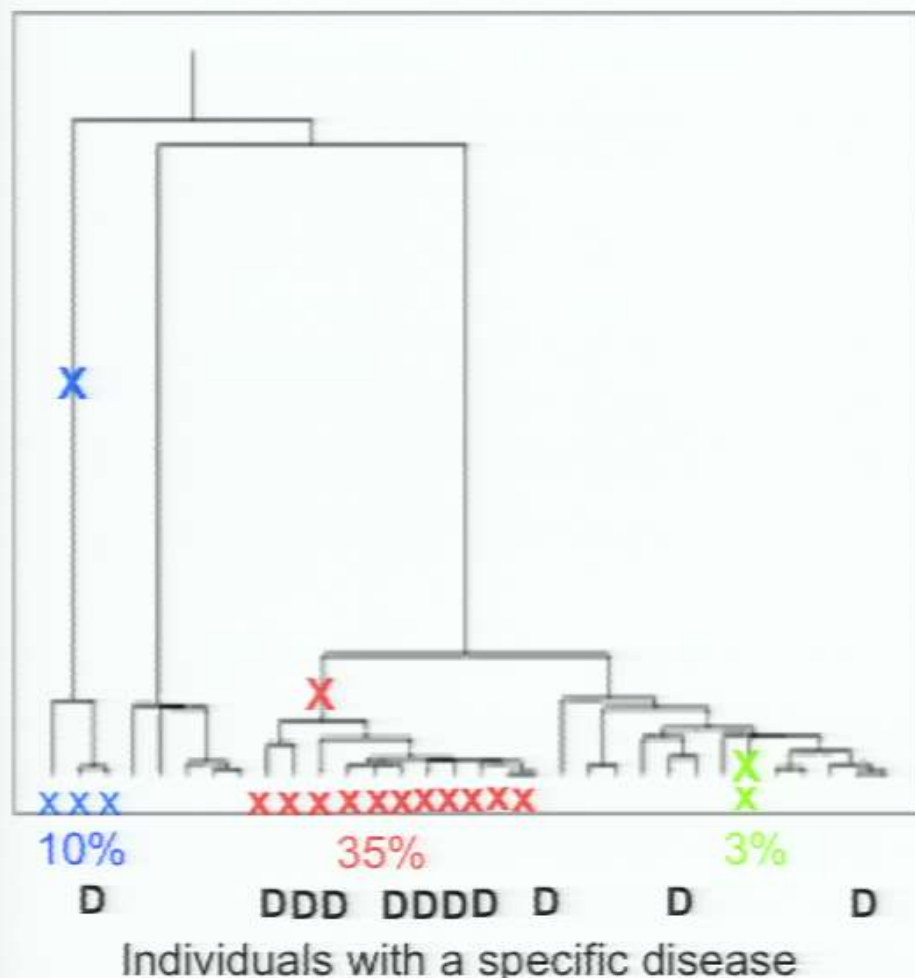


- At each point in the genome we are descended from a common ancestor
- Mutations since the common ancestor give rise to genetic variants shared by the descendants
- Some mutations are associated with disease: the frequency in disease cases is different from in controls

The tree changes along the genome

- At each locus there is a tree
- Ancestral recombinations change the tree as you move along the genome
- The resulting *Ancestral Recombination Graph* describes the way that individual sequences in a sample are related

DNA sequences are related by evolution



- At each point in the genome we are descended from a common ancestor
- Mutations since the common ancestor give rise to genetic variants shared by the descendants
- Some mutations are associated with disease: the frequency in disease cases is different from in controls

The tree changes along the genome

- At each locus there is a tree
- Ancestral recombinations change the tree as you move along the genome
- The resulting *Ancestral Recombination Graph* describes the way that individual sequences in a sample are related

The tree changes along the genome

- At each locus there is a tree
- Ancestral recombinations change the tree as you move along the genome
- The resulting *Ancestral Recombination Graph* describes the way that individual sequences in a sample are related

a ..C..G..A..

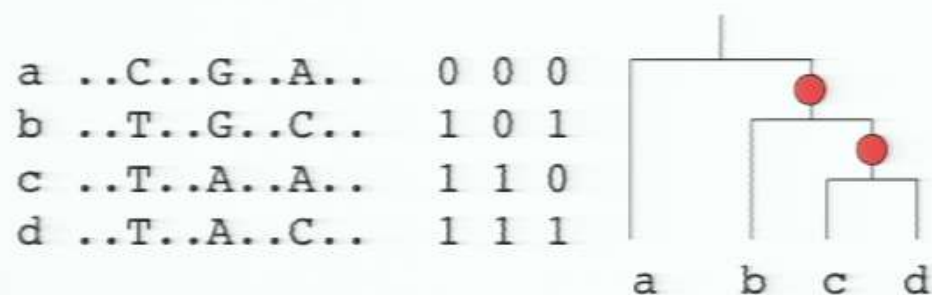
b ..T..G..C..

c ..T..A..A..

d ..T..A..C..

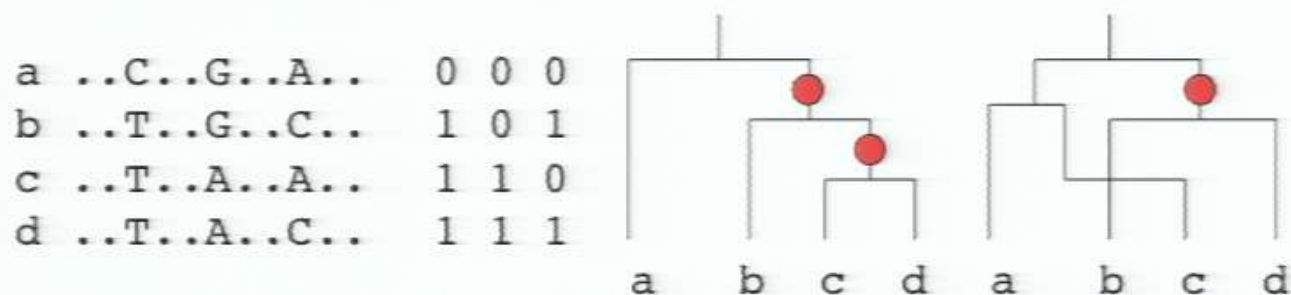
The tree changes along the genome

- At each locus there is a tree
- Ancestral recombinations change the tree as you move along the genome
- The resulting *Ancestral Recombination Graph* describes the way that individual sequences in a sample are related



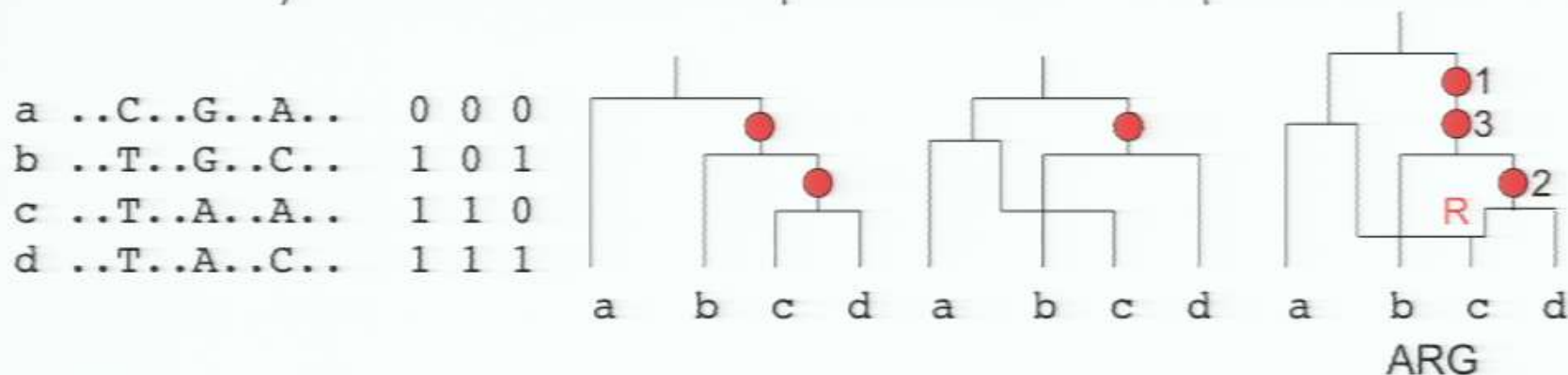
The tree changes along the genome

- At each locus there is a tree
- Ancestral recombinations change the tree as you move along the genome
- The resulting *Ancestral Recombination Graph* describes the way that individual sequences in a sample are related



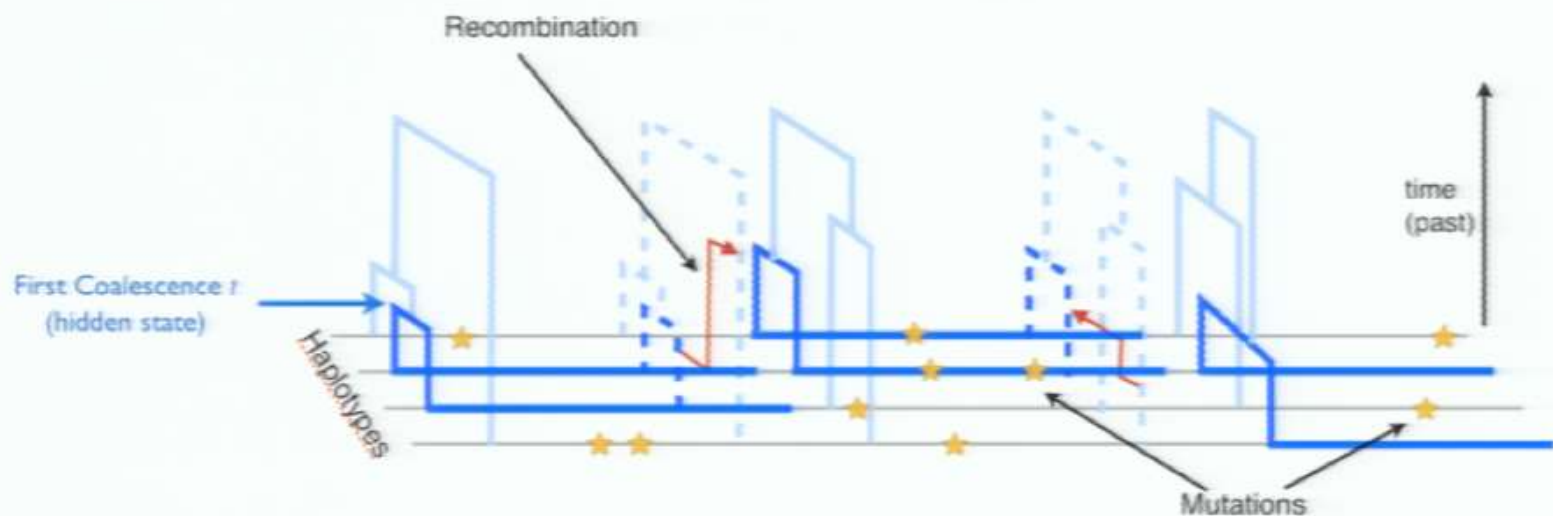
The tree changes along the genome

- At each locus there is a tree
- Ancestral recombinations change the tree as you move along the genome
- The resulting *Ancestral Recombination Graph* describes the way that individual sequences in a sample are related



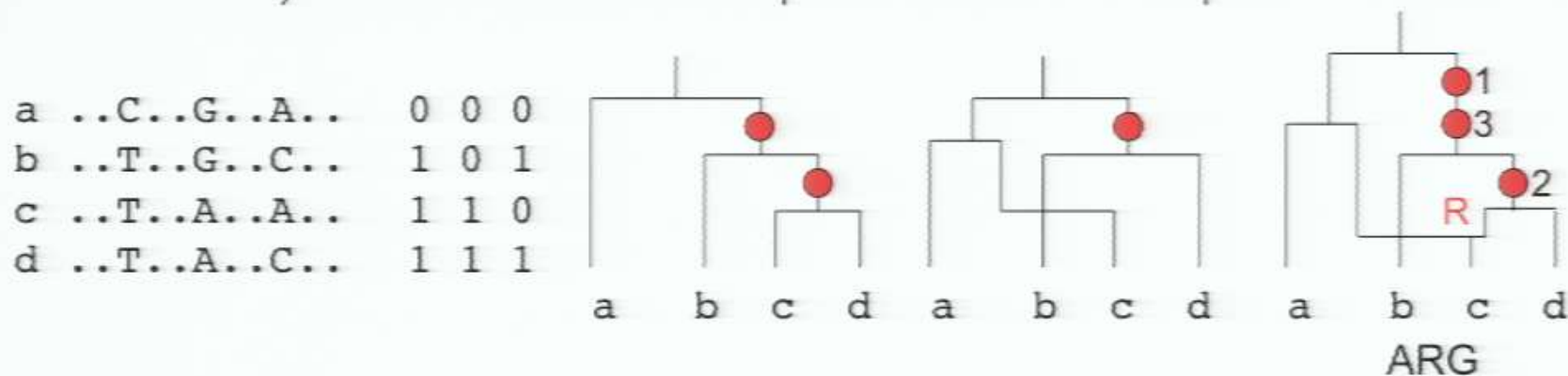
Inference on ARGs is hard

- There is a natural generative model for ARGs
 - Coalescent with recombination (Kingman, Tavaré, Griffiths)
 - “Prune and graft” Markovian approximation is very good
- But infinitely many ARGs satisfy a data set
 - Sampling histories is hard, as is conditional sampling



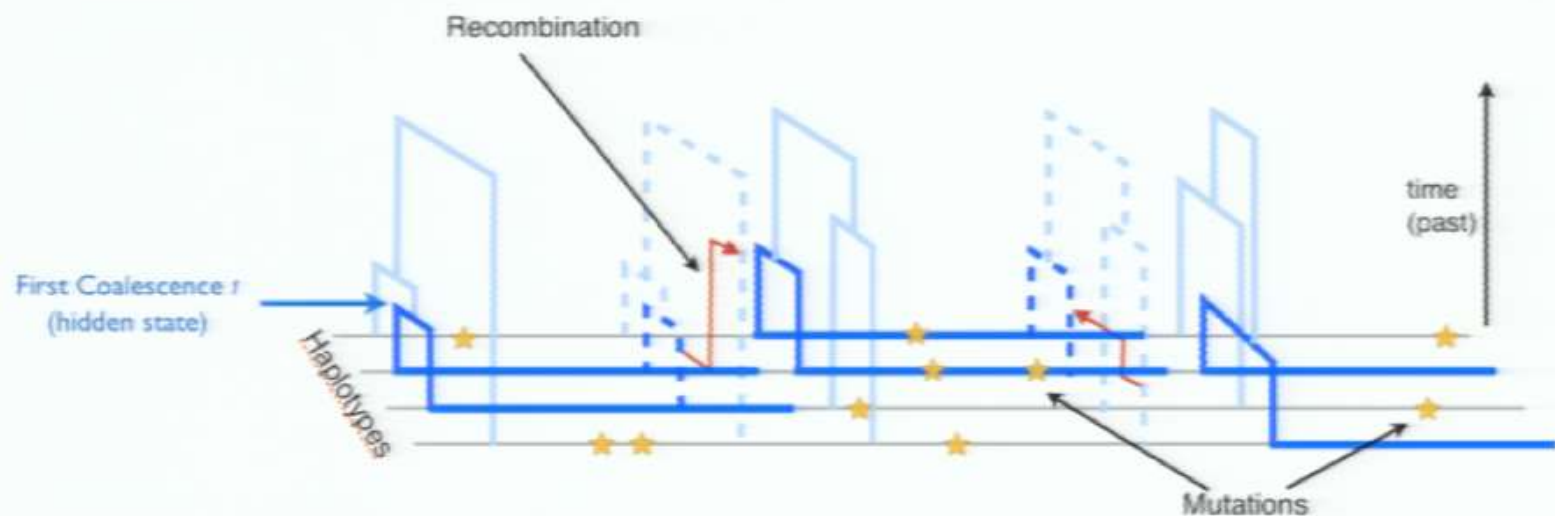
The tree changes along the genome

- At each locus there is a tree
- Ancestral recombinations change the tree as you move along the genome
- The resulting *Ancestral Recombination Graph* describes the way that individual sequences in a sample are related



Inference on ARGs is hard

- There is a natural generative model for ARGs
 - Coalescent with recombination (Kingman, Tavaré, Griffiths)
 - “Prune and graft” Markovian approximation is very good
- But infinitely many ARGs satisfy a data set
 - Sampling histories is hard, as is conditional sampling



Sequencing accuracy requirements

refNCBI36	8881	CCTCCGCGCTCCAGGTTCAACTGATTCTTCTGCCTCAGGCTCCCAAGTAGCTGGGATTACAGGCGAGCGCCACCACACCTGGGCTAAATTTTTGTATTTTTTAGTAGAGACGGGGTTTTG	9008
Venter	8881	9008
Watson	8881	9008
refNCBI36	9001	CCACGTTGGCCAGGCTGGTCTTGAACCTCTGACCTCAGGTGATCTGCCCGGCTTGGGCTCCCAAAGTGCTGGGATTACAGGTGTGAGCTACCGCGCTTAGCCAGTGATAGAGTTTTTGT	9128
Venter	9001	9128
Watson	9001	9128
refNCBI36	9121	TGCCAAAACAAAACATATGAACATATGATAGCTCTAATAAAAAATGCTGTTTCTTTGTTCTCATAATTTCACTAGCTGAACATATGCTCCATTTTCATCTGTAAAAGAGAAATAATCTGTAC	9248
Venter	9121	9248
Watson	9121	9248
refNCBI36	9241	CTTCITGAGTG	9368
Venter	9241	9368
Watson	9241	9368
refNCBI36	9361	ATACTTTATTATTTTAAAGAAATGAGGCACTTGTGATTTTTTCCCCAACATTTTCATCAGAGAGCAGAAAAGGTGCTGCCAAATTCAGACTTCCACATGAAGATTTCATATGCCAGG	9488
Venter	9361	9488
Watson	9361	9488
refNCBI36	9481	TGACCTGCACTAGAAGCAGGTTAATAACCTGTTTTTCATGGGCTGTCTTGGCTTTCACAATGAATGGTCTCCTTTTGCAATGAATTTTGAAGTTTGTTTTTATTATTCTATGTAATAA	9608
Venter	9481	9608
Watson	9481	9608
refNCBI36	9601	TTTGGCTACATGTAATTATTAGCTGCTTTAGAGTCATGGACTTGGAGAAATGAAGATTTTCTGGTCCAATATTTCTCTAAGTATGTTTGGTATCTGTACTTAAAGAGGGGCTCTGGG	9728
Venter	9601	9728
Watson	9601	9728
refNCBI36	9721	ATCAATTAGGTCAGGAAACTGCTGAGTTAAAAAGATTTTAGTTTTTGGGGAGTTGGGGAGGTGAGAACATTCCTGAGACTTTAATATGCTAATATGAGTCTCTCAGAGAAACAAACA	9848
Venter	9721	9848
Watson	9721	9848
refNCBI36	9841	TACAGTGTTTTACAAACTTACTTGACCACAAAACCTTTTCTCGTGGTCTTGTGTACTCAGAACTCAAATCTTTGAATTTACAAATGAACAGAAACGACAGCAATGGCTTATTCAA	9968
Venter	9841	9968
Watson	9841	9968
refNCBI36	9961	AGCAATATAGCAGTTTATTAAAGAACTGGACTTGAACATA	10008
Venter	9961	10008
Watson	9961	10008

Sequencing accuracy requirements

refNCBI36	8881	CCTCCGCTCCAGGTTCAACTGATTCTTCTGCTCAGCTCCCAAGTAGCTGGGATTACAGGCGAGCGCCACCACACCTGGCTAATTTTTTGTATTTTTTAGTAGAGACGGGGTTTTG	9008
Venter	8881	9008
Watson	8881	9008
refNCBI36	9001	CCACGTTGGCCAGGCTGGTCTTGAACCTCTGACCTCAGGTGATCTGCCCGCTTGGCCTCCCAAGTGCTGGGATTACAGGTGTGAGCTACCGCGCTTAGCCAGTGATAGAGTTTTTGT	9128
Venter	9001	9128
Watson	9001	9128
refNCBI36	9121	TGCCAAACAAACATATGAACATATGATAGCTCTAATAAAAAATGCTGTTTCTTTGTTCTCATAATTCAGTAGCTGAACATATGCTCCATTTTCATCTGTAAAGAGAAATAATCTGTAC	9248
Venter	9121	9248
Watson	9121	9248
refNCBI36	9241	CTTCTTGAGTGGAGAGGGAGAAAGGAATGAAAGCAAAATAACTATTTATANTTAGTGAAGAGAAAGCTTAGGAAAAATGAGAGAGTCTTTAGAAATGTTAAATAGCTTTAATAATA	9368
Venter	9241	9368
Watson	9241	9368
refNCBI36	9361	ATACTTTATTTATTTAAAGAAATGTAGGCAGTTGTGATTTTTCCCCAACATTTTCATCAGAGAGCAGAAAAGGTGCTGCCAAATTCAGACTTCCACATGAAGATTTCACTATGCCAGG	9488
Venter	9361	9488
Watson	9361	9488
refNCBI36	9481	TGACCTGCACT	9608
Venter	9481	9608
Watson	9481	9608
refNCBI36	9601	TTTGGCTACATGTAATTATTAGCTGCTTTACAGTCATGGACTTGGAGAAATCAAGATTTTCTGGTCCAATATTTCTAAGTATGTTCTGTCATCTGTACTTAAAGAGGGGCTCTGGG	9728
Venter	9601	9728
Watson	9601	9728
refNCBI36	9721	ATCAATTAGGTCAGGAAACTGCTGAGTTAAAAAGATTTTAGTTTTTGGGGAGTTGGGGAGGTGAGAACATTCCTGAGACTTTAATATGCTAATATGAGTCTCTCAGAGAAACACA	9848
Venter	9721	9848
Watson	9721	9848
refNCBI36	9841	TACAGTGTTTTACAAACTTACTTGACCACAAACCITTTTCTCGTGGTCTTGTGTTACTCAGAACTCAAATCTTTGAATTTACAAATGAACAGAAACGACAGCAANTGGCTTATTCAA	9968
Venter	9841	9968
Watson	9841	9968
refNCBI36	9961	AGCAATATAGCAGTTTATTTAAAGAACTGGACTTGAACATA	10008
Venter	9961	10008
Watson	9961	10008


1% false positives implies 10^{-5} errors per bp

Sequencing accuracy requirements

refNCBI36	8881	CCTCCGCGCTCCAGGTTCAACTGATTCCTCTGCGCTCAGGCTCCCAAGTAGCTGGGATTACAGGCGAGCGCCACCACACCTGGCTAAATTTTTTGTATTTTTTTAGTAGAGACGGGGTTTTG	9008
Venter	8881	9008
Watson	8881	9008
refNCBI36	9001	CCACGTTGGCCAGGCTGGTCTTGAAGTCTGACCTCAGGTGATCTGCCCGCCTTGGCCTCCCAAAGTGCTGGGATTACAGGTGTGAGCTACCGCGCTTAGCCAGTGATAGAGTTTTTGT	9128
Venter	9001	9128
Watson	9001	9128
refNCBI36	9121	TGCCAAACAAACATATGAACATATGATAGCTCTAATAAAAAATGCTGTTCCTTTGTTCTCATAATTCAGTAGCTGAAGTATGCTCCATTTTCATCTGTAAAAGAGAAATAATCTGTAC	9248
Venter	9121	9248
Watson	9121	9248
refNCBI36	9241	CTTCCTGAGTGGAGAGGGAGAAAGGAATGAAAGCAAAATTAATATTTATTAATTAGTGAAGAGAAAGCTTAGGAAAAAATGAGAGAGTGCCTTTAGAAATGTTAAATAGCTTTAATAATA	9368
Venter	9241	9368
Watson	9241	9368
refNCBI36	9361	ATACTTTATTATTTTAAAGAAATGTAAGCAGTTGTGATTTTTTCCCCAACATTTTCATCAGAGAGCAGAAAGGTGCTGCCAAATTCAGACTTCCACATGAAGATTTCAGTATGCCAGG	9488
Venter	9361	9488
Watson	9361	9488
refNCBI36	9481	TGACCTGCACT	9608
Venter	9481	9608
Watson	9481	9608
refNCBI36	9601	TTTGGCTACATGTAATTATTAGCTGCCCTTACAGTCATGCACTTGGAGCAATTGAAGATTTTCTGGTCCAATATTTCTCTAAGTATGTTCTGTCATCTGTACTTAAAGAGGGGCTCTGGG	9728
Venter	9601	9728
Watson	9601	...	9728
refNCBI36	9721	AT	9848
Venter	9721	...	9848
Watson	9721	9848
refNCBI36	9841	TACAGTGTTTTACAAACTTACTTGACCACAAAACCTTTTCTCGTGGTCTTGTGTACTCAGAACTCAAATCTTTTGAATTTACAAATGAACAGAAACGACAGCAATGGCTTATTCAA	9968
Venter	9841	9968
Watson	9841	9968
refNCBI36	9961	AGCAATATAGCAGTTTATTTAAAGAACTGGACTTGAAGTA	10008
Venter	9961	10008
Watson	9961	10008

1% false positives implies 10^{-5} errors per bp

>1000 individuals implies $<10^{-8}$ errors per bp per person



1000 Genomes

A Deep Catalog of Human Genetic Variation

- International project to construct a next generation open access baseline data set for human genetics
 - Consortium with multiple centres, platforms, funders
 - Conceived in 2007, pilot projects published Oct 2010, phase I with 1092 samples published November 2012, **Completion 2014**
- Aims
 - Find >95% accessible SNPs at allele frequencies above 1%, down towards 0.1% in coding regions
 - Also discover and characterize indels, structural variants
 - Identify a reference set of human genome sequences
- Driver for data and methods

1000 Genomes

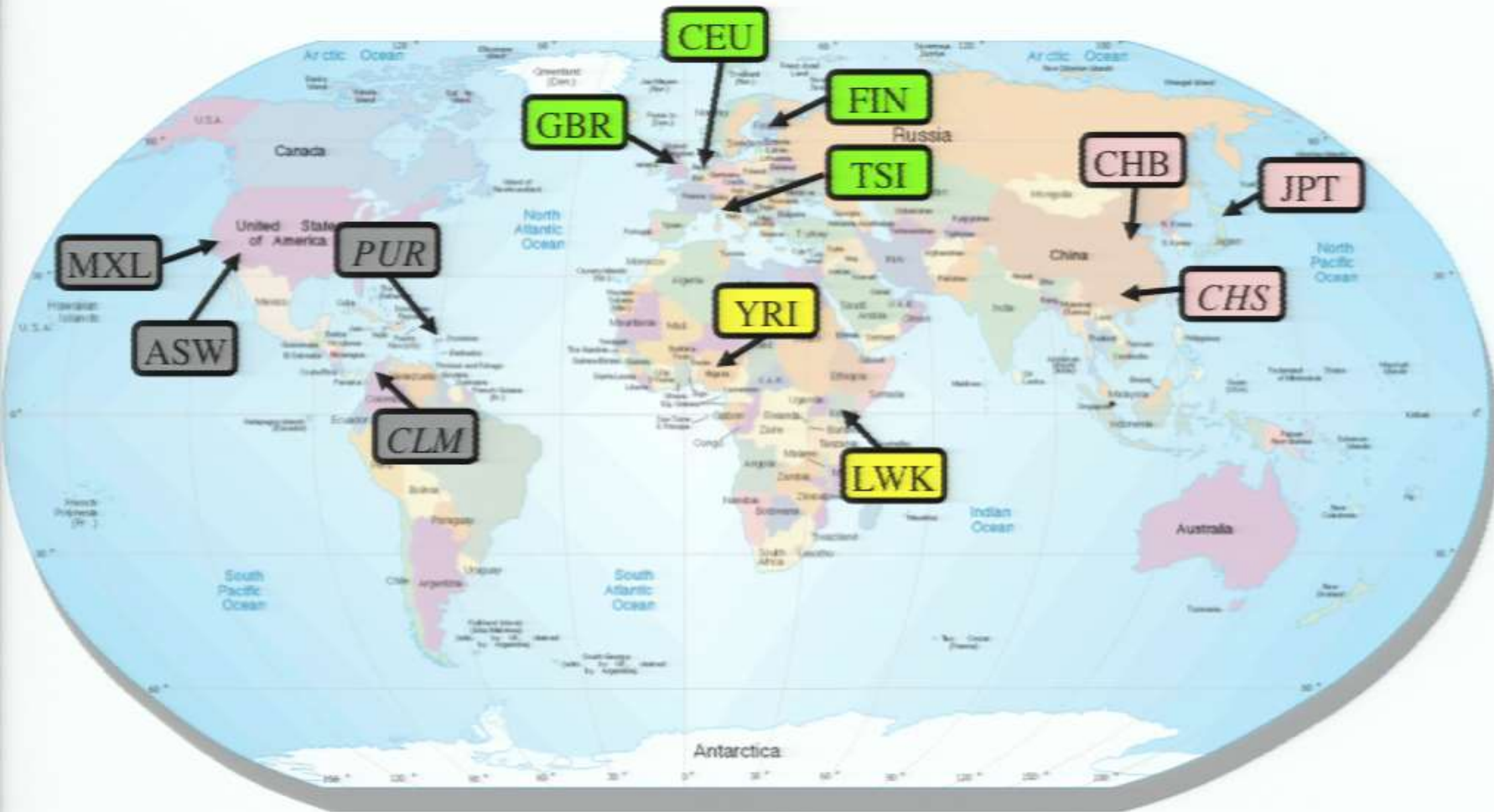
Pilot project 179 samples (Nature 2010)



~100 per population: 4x Whole Genome Shotgun + Deep Exomes

1000 Genomes

Pilot project 179 samples (Nature 2010)
Phase I: 1,092 samples (Nature 2012)



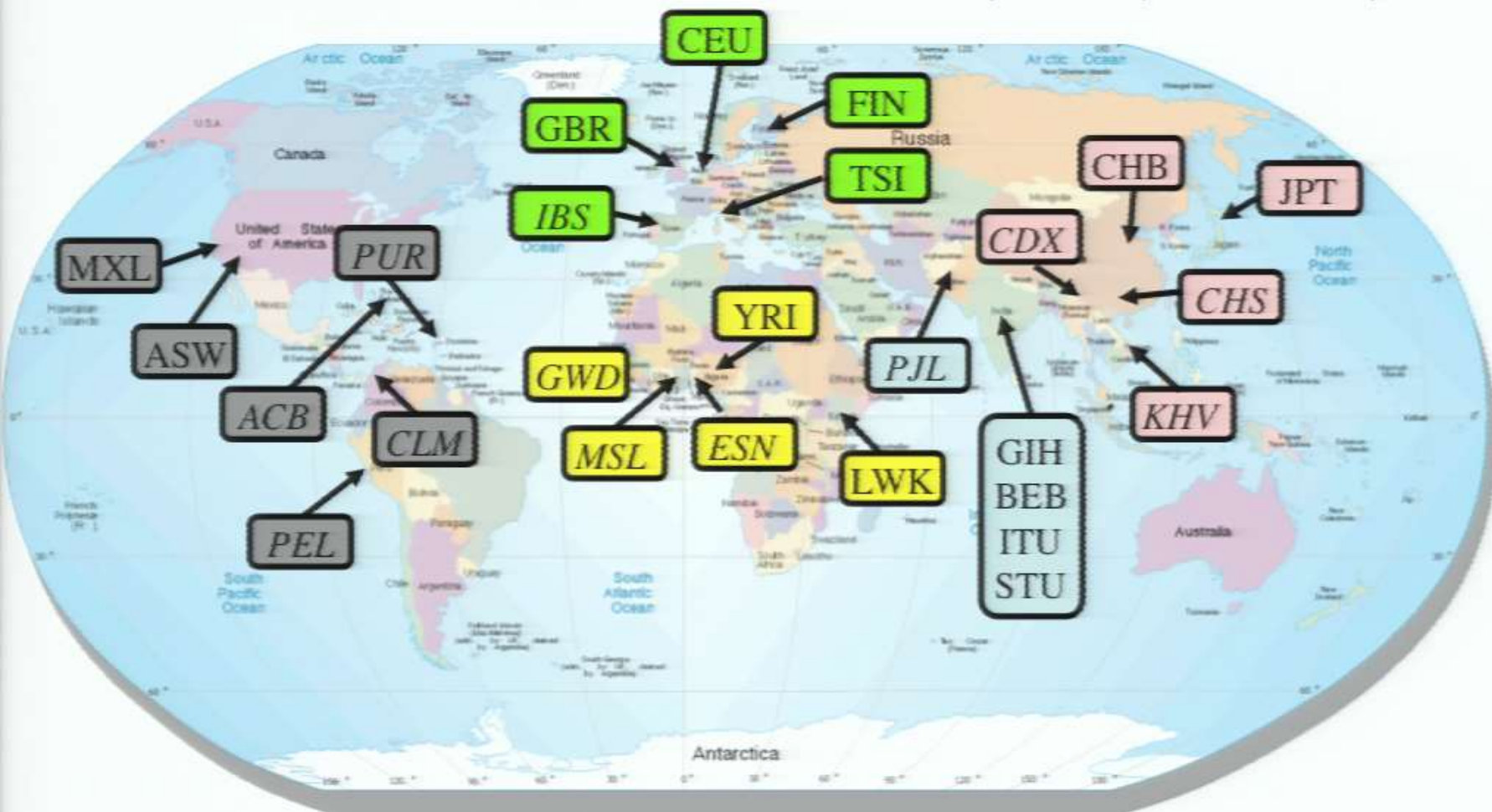
~100 per population: 4x Whole Genome Shotgun + Deep Exomes

1000 Genomes

Pilot project 179 samples (Nature 2010)

Phase 1: 1,092 samples (Nature 2012)

Phase 3: >2,500 samples sequence complete

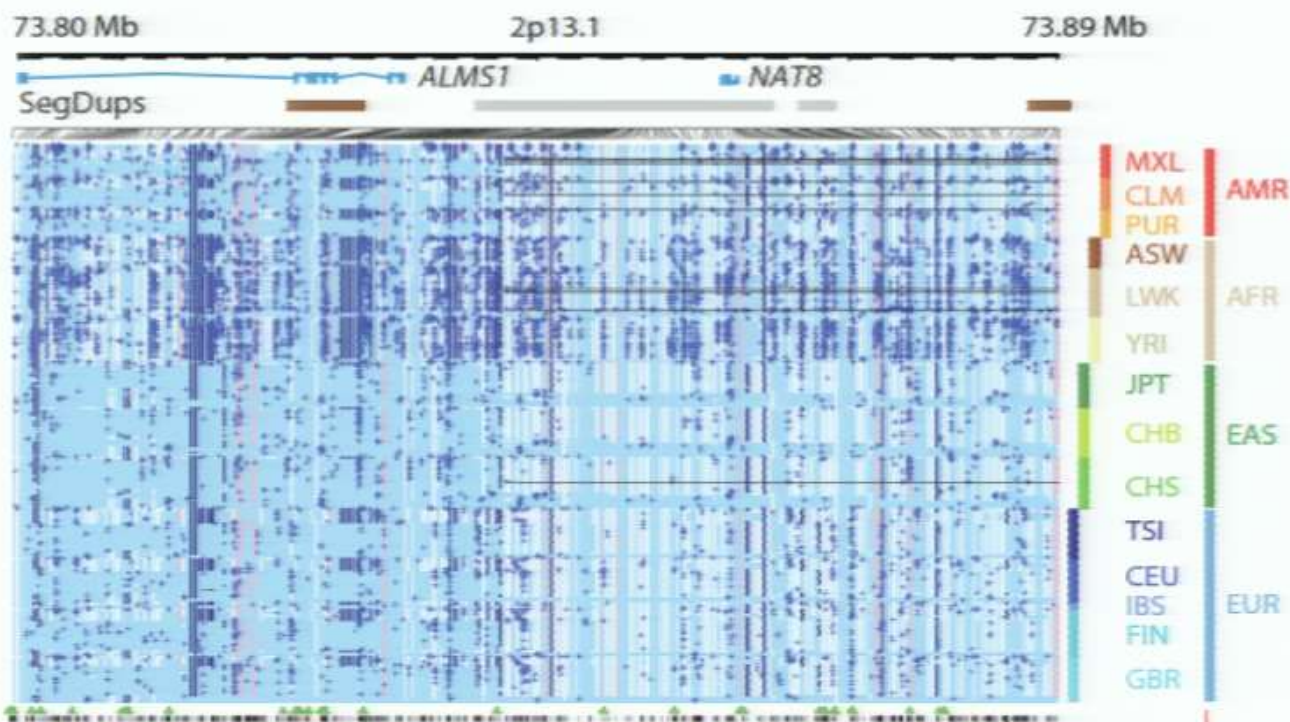


~100 per population: 4x Whole Genome Shotgun + Deep Exomes

1000 Genomes

A wealth of data

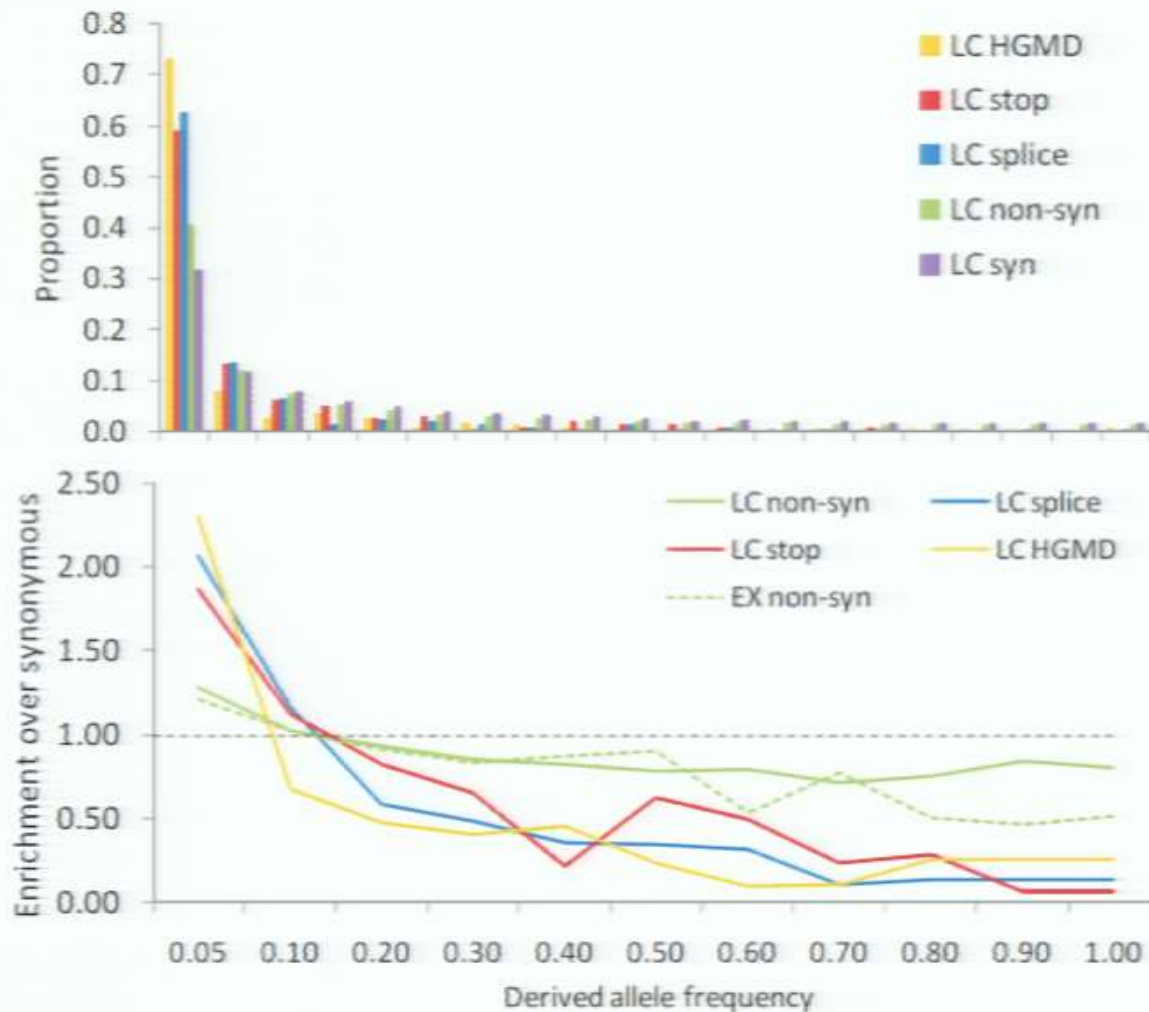
Call set	Samples	SNPs	Sensitivity (HapMap3.3)	FP (in 59,721 monomorphic OMNI sites)	Genotype accuracy (HapMap 3 hets)
Pilot	179	15.2M	97.65%	43,606	97.75%
Phase I	1,092	38.0M	98.87%	1,261	99.24%



Also integrate
 ~500k exome calls,
 1.4M indels,
 14k large deletions,
 Omni2.5M
 genotypes

20.2 Tbp sequence

Distribution of functional variants



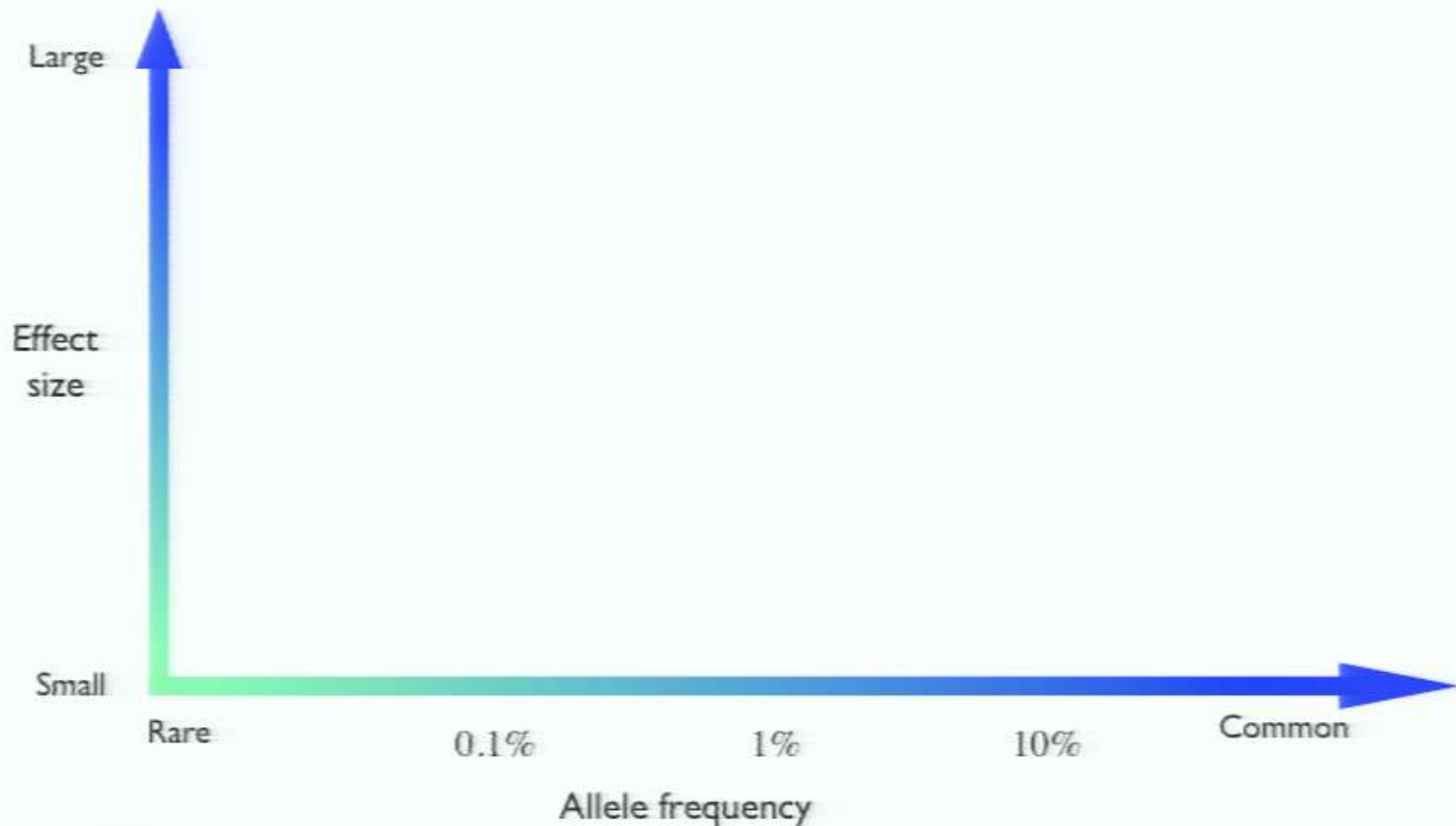
Common deleterious mutations are removed by selection

Pilot paper

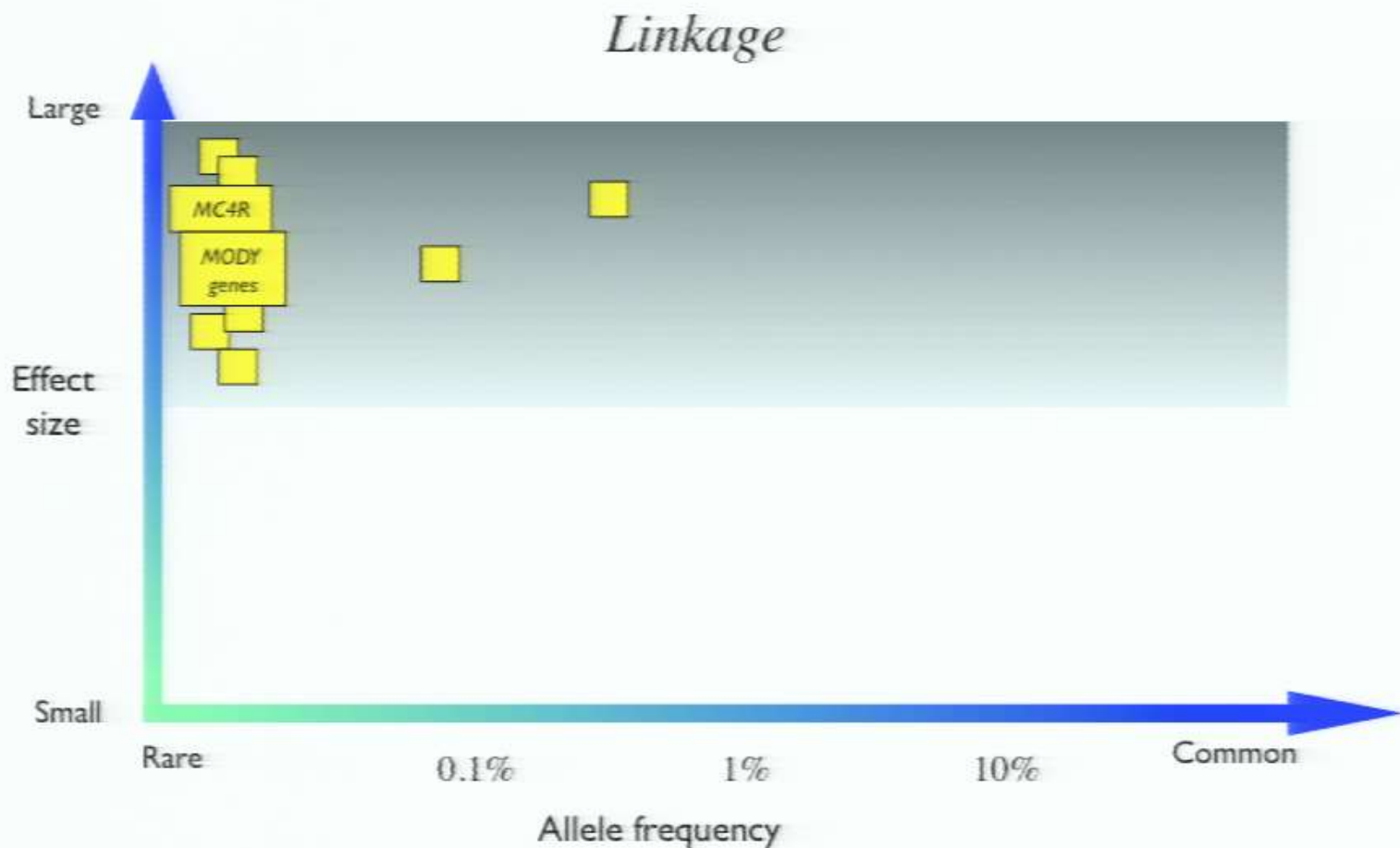
Per individual deleterious variant load

Variant type	Number of variants (pop range)	Excess rare deleterious <0.5%	Excess low freq deleterious 0.5%-5%
All sites	3.75M-4.7M		
Conserved sites (GERP>2)	140k-200k	150-510	250-1.3k
Synonymous conserved	1.4k-1.9k		
NonSyn conserved	2.7k-4k	76-190	77-130
Loss-Of-Function conserved	116-175	6-13	11-31
Non-coding RNA conserved	200-290	1-3	4-13
Motif loss in TF peak (incomplete)	670-1020	8-22	20-110

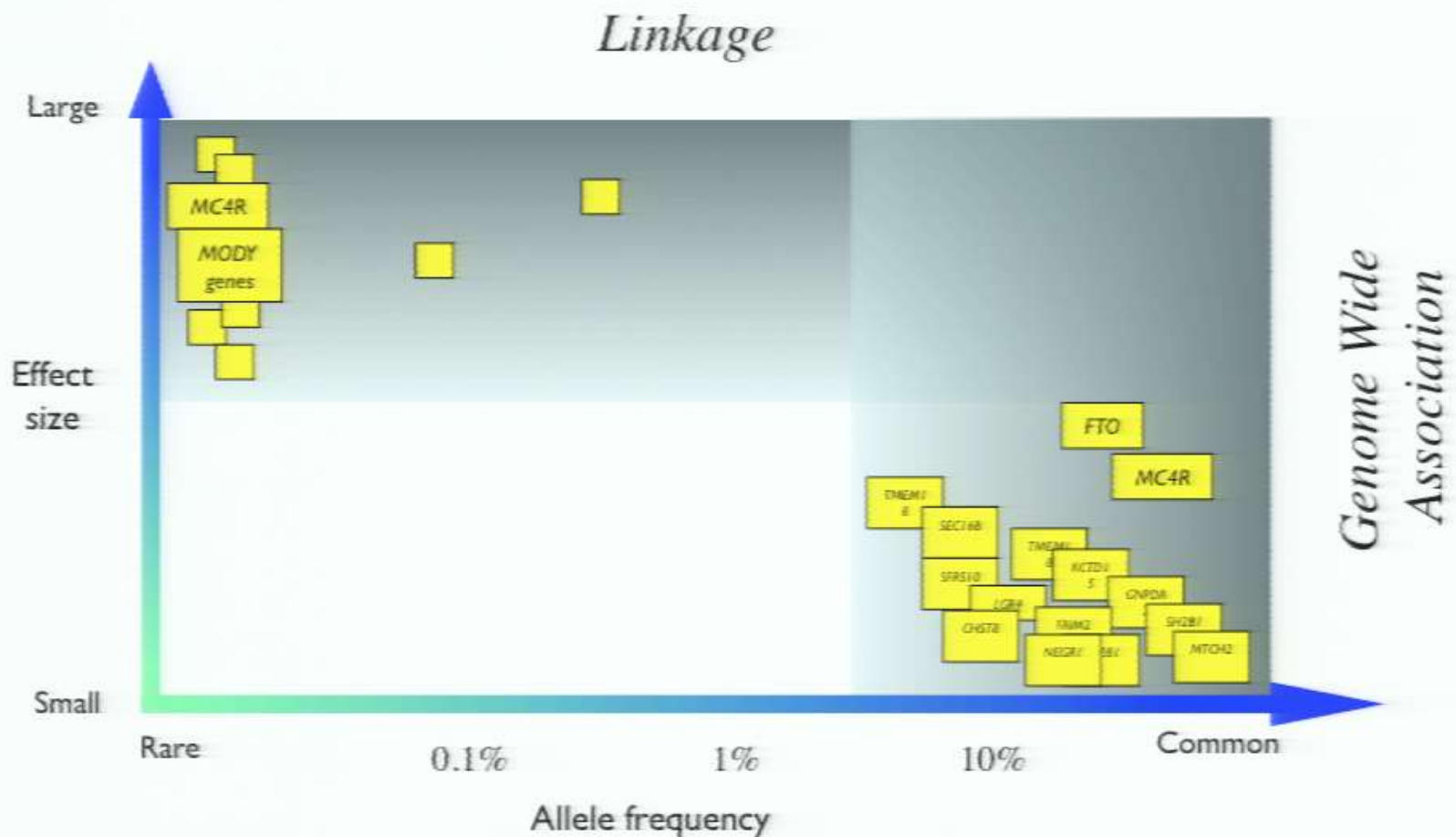
What about identifying functional genetic variants?



What about identifying functional genetic variants?

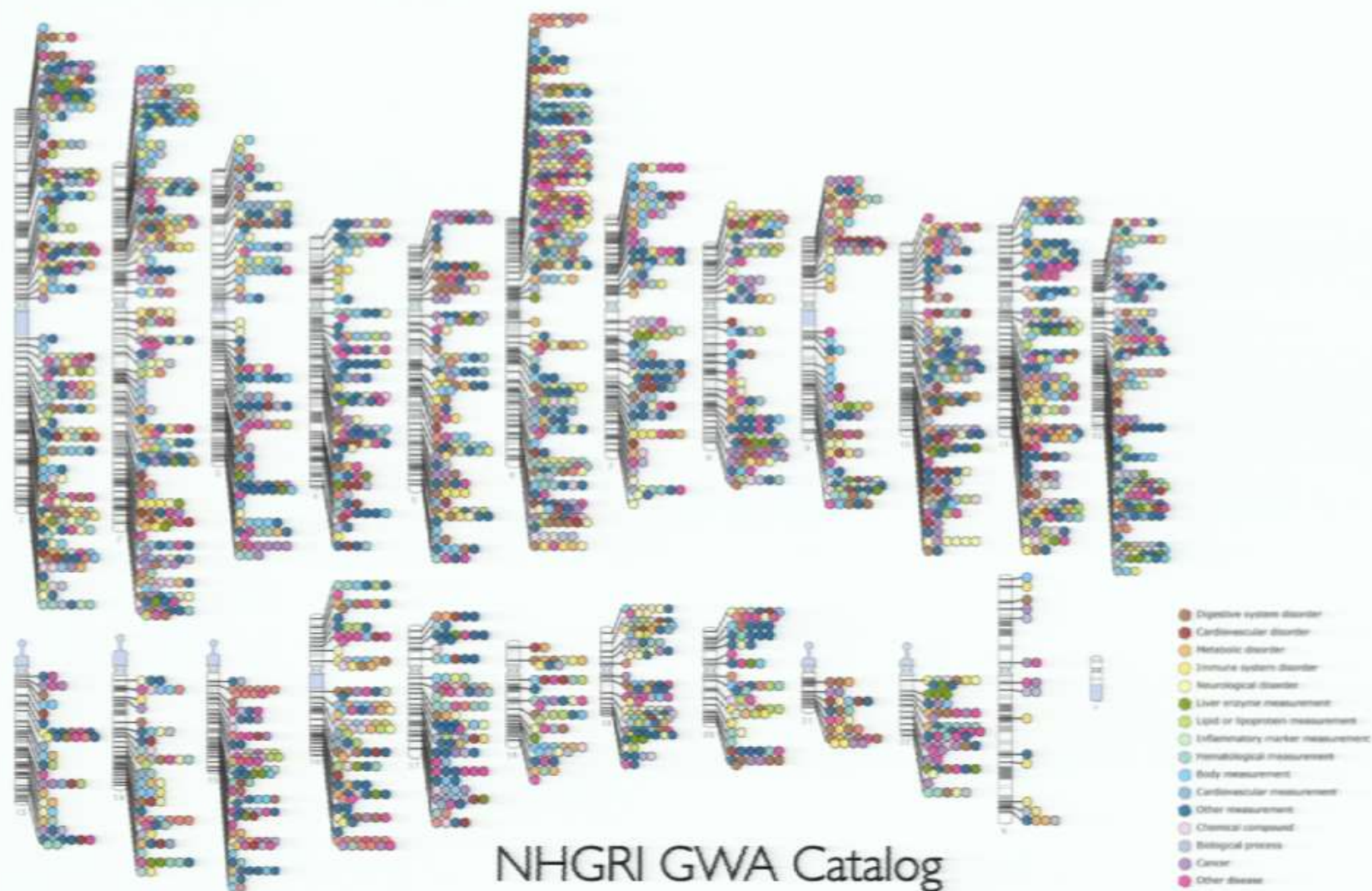


What about identifying functional genetic variants?



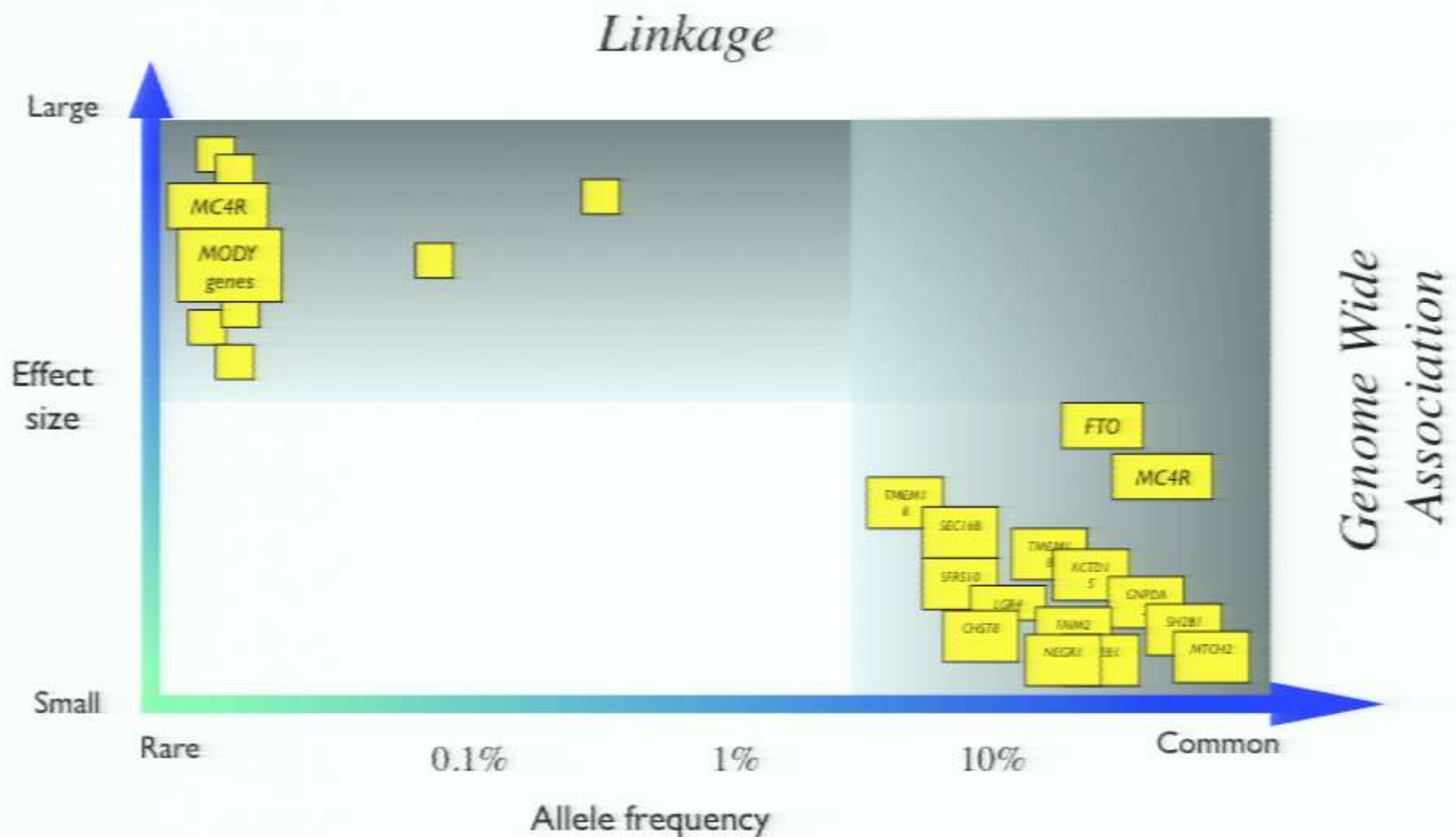
Over a million people genotyped

Published Genome-Wide Associations through 07/2012

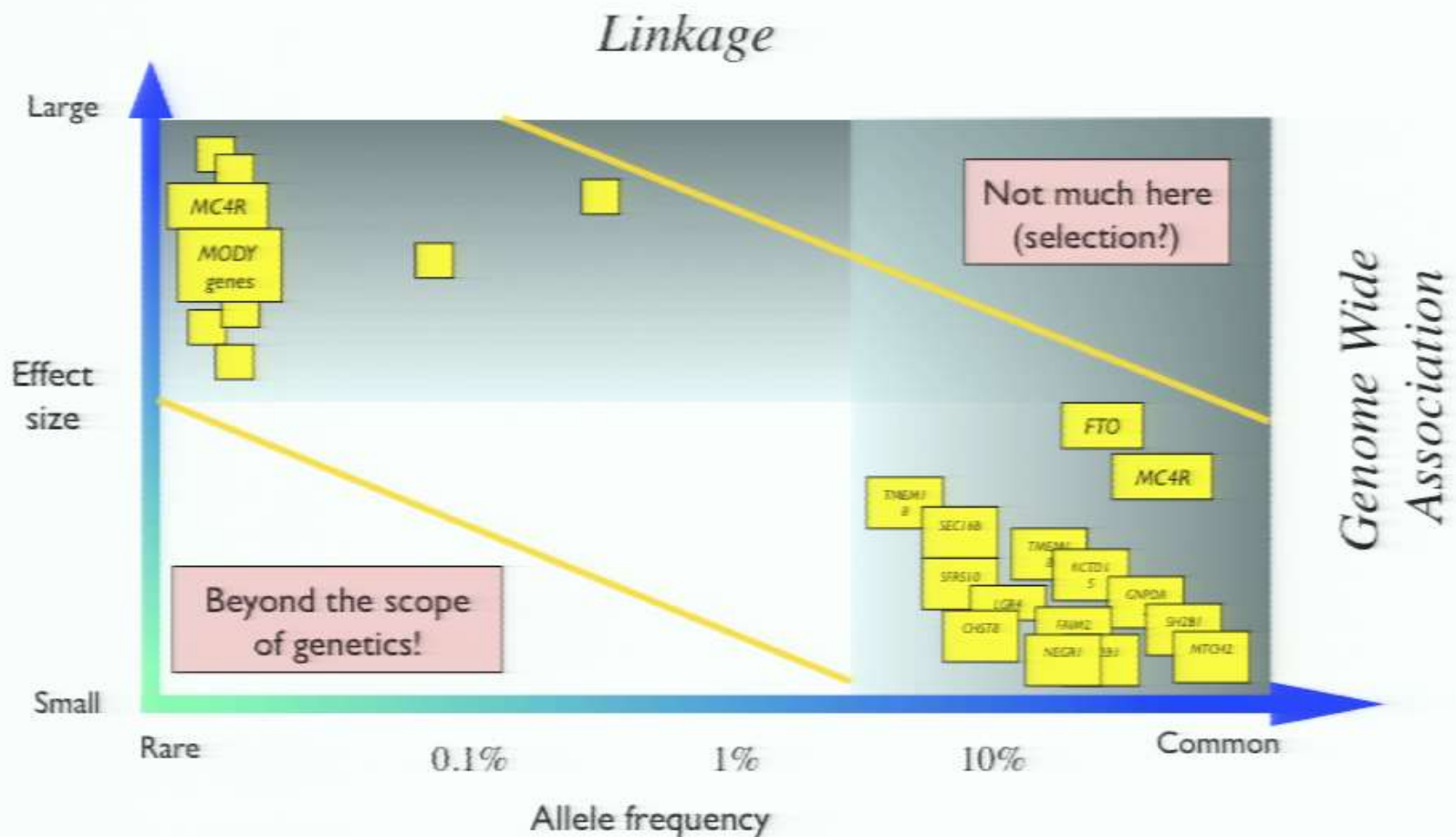


NHGRI GWA Catalog
www.genome.gov/GWAStudies

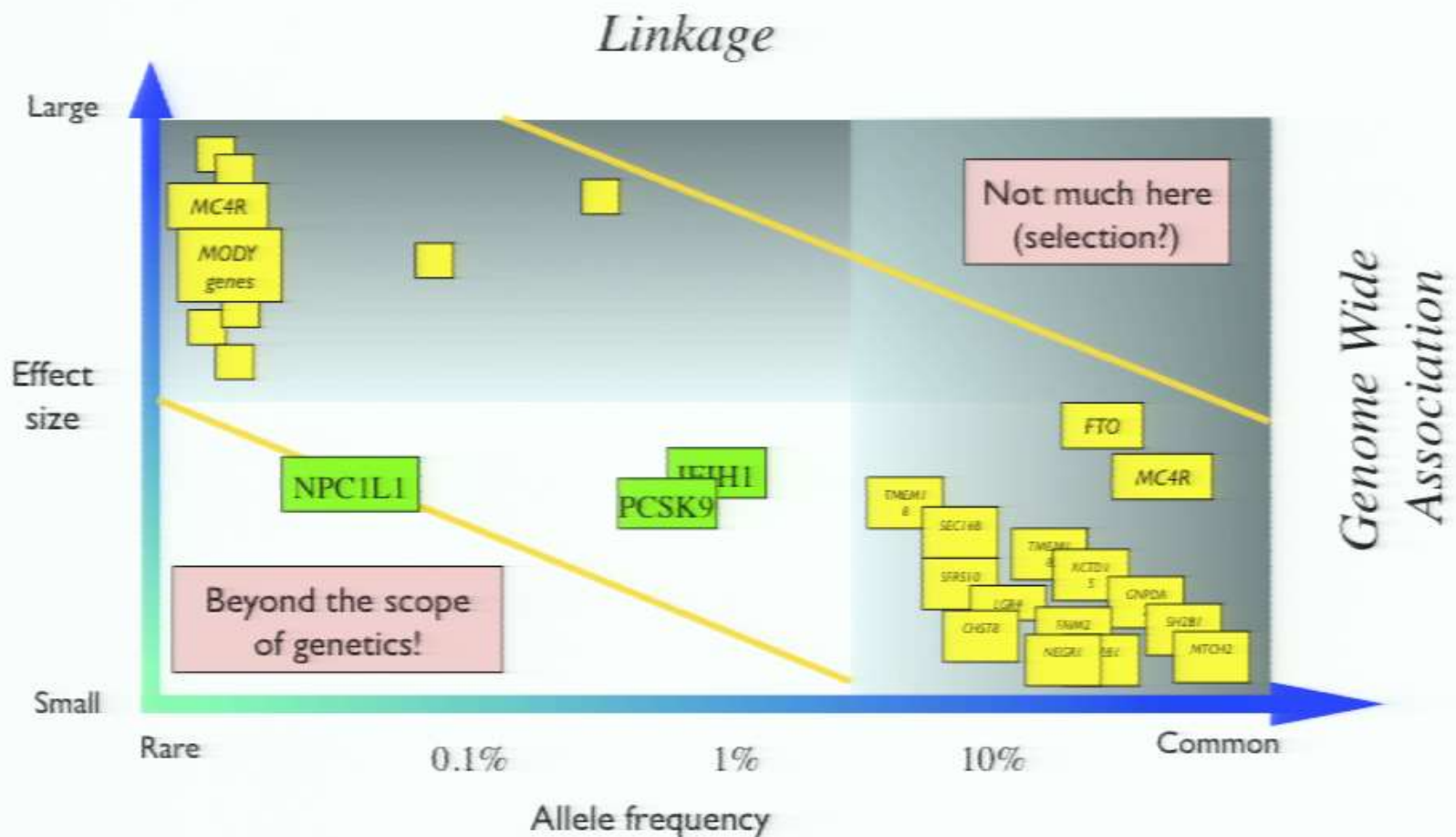
What about identifying functional genetic variants?



What about identifying functional genetic variants?



What about identifying functional genetic variants?



A few genes in the middle zone have been found by candidate studies

10,000 UK Genomes (2010-13)

- Sequence
 - 4,000 cohort samples genome wide: TwinsUK and ALSPAC
 - 6,000 exomes from samples with extreme phenotypes
- Goals
 - Direct association in sequenced samples
 - Impute new variants into additional samples, GWAS sets
 - Provide a sequence variation resource for use in further studies
- Progress
 - Data collection started late 2010
 - 4,004 WGS samples sequenced, ~6,300 exomes sequenced
 - Calls on ~3,700 WGS and ~6,200 exomes available
 - Greater than 50M variants, many novel

Cohorts initial analysis: 54 traits in 14 groups on 2,453 samples

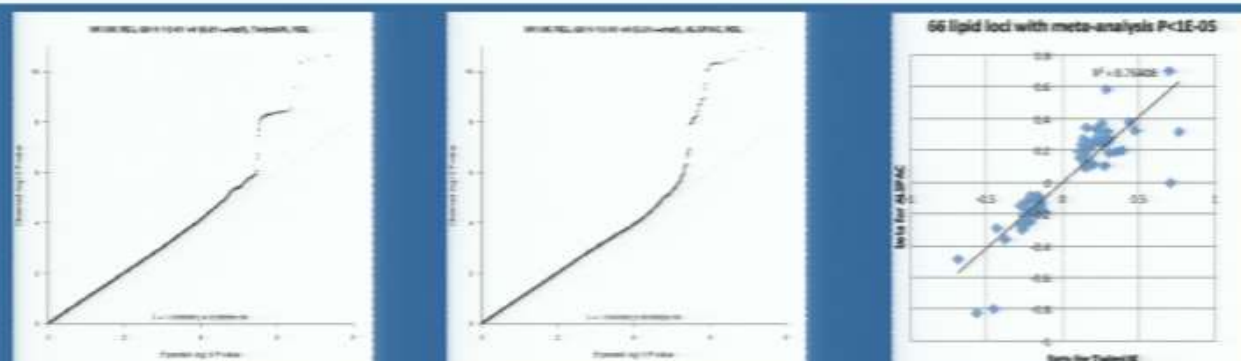


Fig 3: QQplots and effect size comparison of singlepoint analysis of HDL (MAF>1%)

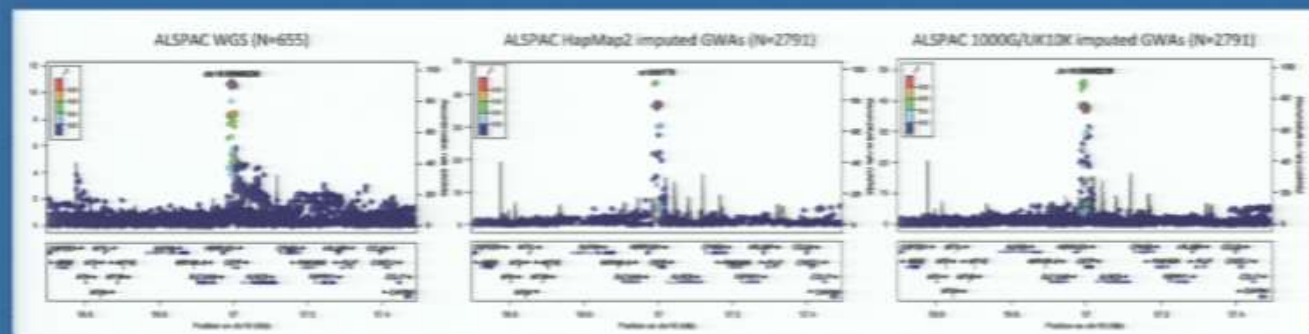


Fig 4: Association between *CETP* SNPs and HDL for ALSPAC WGS (N=655), ALSPAC HapMap2 imputed GWAs (N=2,791) and ALSPAC 1000G/UK10K imputed GWAs (N=2,791).

Some potential candidate new loci < 2.5% or < 10% poorly tagged

Alternative strategy: Impute full sequences into GWAS studies from reference sequences

a g a g t t g a g g g a a c c t g a g a a
t g a g a c g a g g g a a a t t g a g a c
t g c g a c g g t g a t t c t c c a g a c
a g c g a c g a t g g t a c t t g a t c a
t a a g t t a g t a a t t c c c g a g c a
t g c a a t g a g g g a a a t t g t t a a
a g a g a c g g g g g a a a t t c t g c c

**Reference haplotypes
via sequencing
studies**

e.g. 1000 GP, UK10K

g a g g t a a
g c t a t t c

a t g g t t a
g c g g c a a

g c t g t t c
g c g g c a a

g t g g t a c
a t t a c a a

**Genotype data from
GWAS**

e.g. T2D, Arthritis, ...

Imputing missing data from reference sequences

a	g	a	g	t	t	g	a	g	g	g	a	a	c	c	t	g	a	g	a	a
t	g	a	g	a	c	g	a	g	g	g	a	a	a	t	t	g	a	g	a	c
t	g	c	g	a	c	g	g	t	g	a	t	t	c	t	c	c	a	g	a	c
a	g	c	g	a	c	g	a	t	g	g	t	a	c	t	t	g	a	t	c	a
t	a	a	g	t	t	a	g	t	a	a	t	t	c	c	c	g	a	g	c	a
t	g	c	a	a	t	g	a	g	g	g	a	a	a	t	t	g	t	t	a	a
a	g	a	g	a	c	g	g	g	g	g	a	a	a	t	t	c	t	g	c	c

**Reference haplotypes
via sequencing
studies**

eg. 1000 Genomes Project

g	a	g	g	t	a	a
g	c	t	a	t	t	c

a	t	g	g	t	t	a
g	c	g	g	c	a	a

g	c	t	g	t	t	c
g	c	g	g	c	a	a

g	t	g	g	t	a	c
a	t	t	a	c	a	a

Imputation of unobserved alleles via matching of shared haplotypes

Imputing missing data from reference sequences

a	g	a	g	t	t	g	a	g	g	g	a	a	c	c	t	g	a	g	a	a
t	g	a	g	a	c	g	a	g	g	g	a	a	a	t	t	g	a	g	a	c
t	g	c	g	a	c	g	g	t	g	a	t	t	c	t	c	c	a	g	a	c
a	g	c	g	a	c	g	a	t	g	g	t	a	c	t	t	g	a	t	c	a
t	a	a	g	t	t	a	g	t	a	a	t	t	c	c	c	g	a	g	c	a
t	g	c	a	a	t	g	a	g	g	g	a	a	a	t	t	g	t	t	a	a
a	g	a	g	a	c	g	g	g	g	g	a	a	a	t	t	c	t	g	c	c

**Reference haplotypes
via sequencing
studies**

eg. 1000 Genomes Project

a	g	a	g	t	a	g	a	g	g	t	a	c	t	t	g	a	t	c	a	
t	g	c	g	a	c	g	g	t	g	a	t	t	c	t	t	c	t	g	c	c

Hidden Markov Model:
"select and copy",
e.g. IMPUTE, MACH

t	a	a	a	a	t	g	a	g	g	g	a	a	a	t	t	g	t	t	a	a
t	g	a	g	a	c	g	a	g	g	g	a	a	c	c	c	g	a	g	c	a

a	g	c	g	a	c	g	a	t	g	g	t	a	a	t	t	c	t	g	c	c
a	g	a	g	a	c	g	a	g	g	g	a	a	c	c	t	g	a	g	a	a

Effectively an
approximation to
conditional sampling
on the ARG

t	g	c	a	a	t	g	a	g	g	g	a	a	a	t	t	g	a	g	a	c
t	a	a	g	t	t	a	g	t	a	a	t	t	c	c	t	g	a	t	c	a

Imputation of unobserved alleles via matching of shared haplotypes

Imputing missing data from reference sequences

a g a g t t g a g g g a a c c t g a g a a
t g a g a c g a g g g a a a t t g a g a c
t g c g a c g g t g a t t c t c c a g a c
a g c g a c g a t g g t a c t t g a t c a
t a a g t t a g t a a t t c c c g a g c a
t g c a a t g a g g g a a a t t g t t a a
a g a g a c g g g g g a a a t t c t g c c

a g a g t a g a g g g t a c t t g a t c a
t g c g a c g g t g a t t c t t c t g c c

t a a a a t g a g g g a a a t t g t t a a
t g a g a c g a g g g a a c c c g a g c a

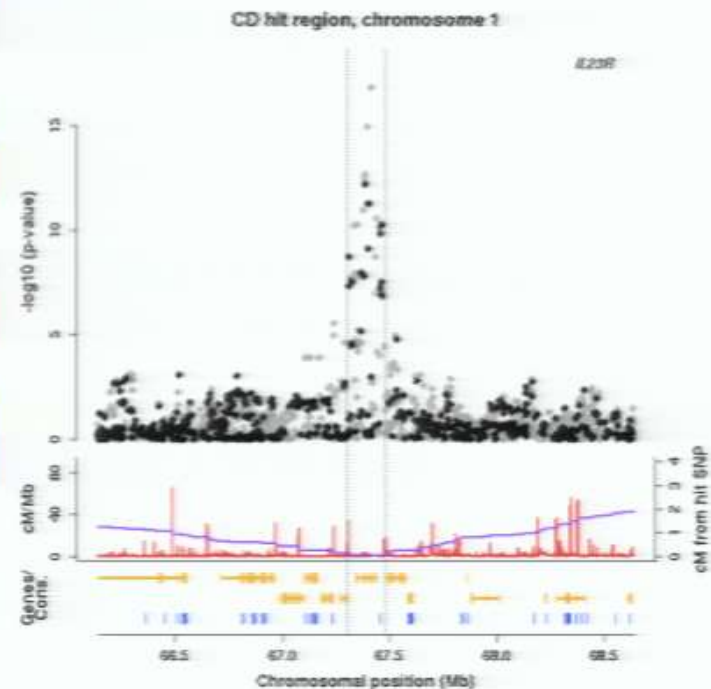
a g c g a c g a t g g t a a t t c t g c c
a g a g a c g a g g g a a c c t g a g a a

t g c a a t g a g g g a a a t t g a g a c
t a a g t t a g t a a t t c c t g a t c a



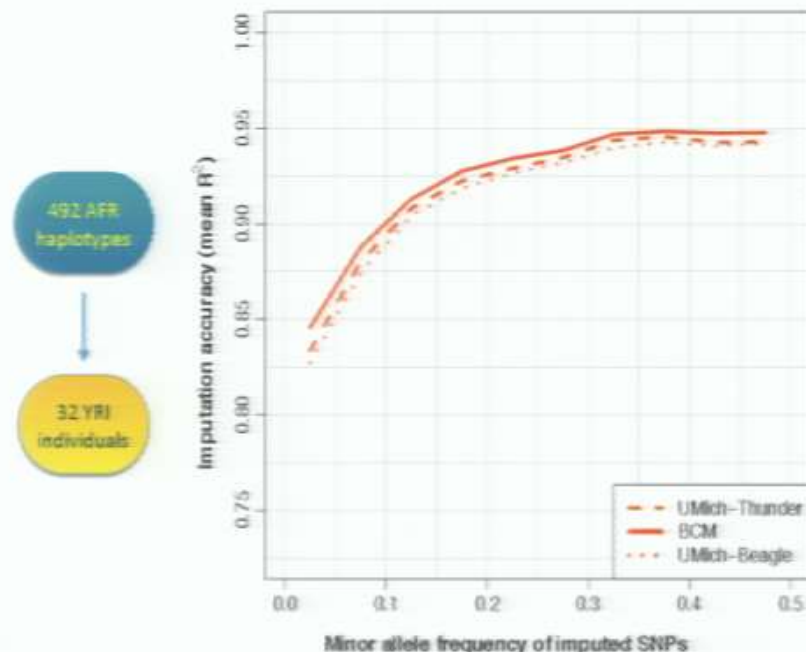
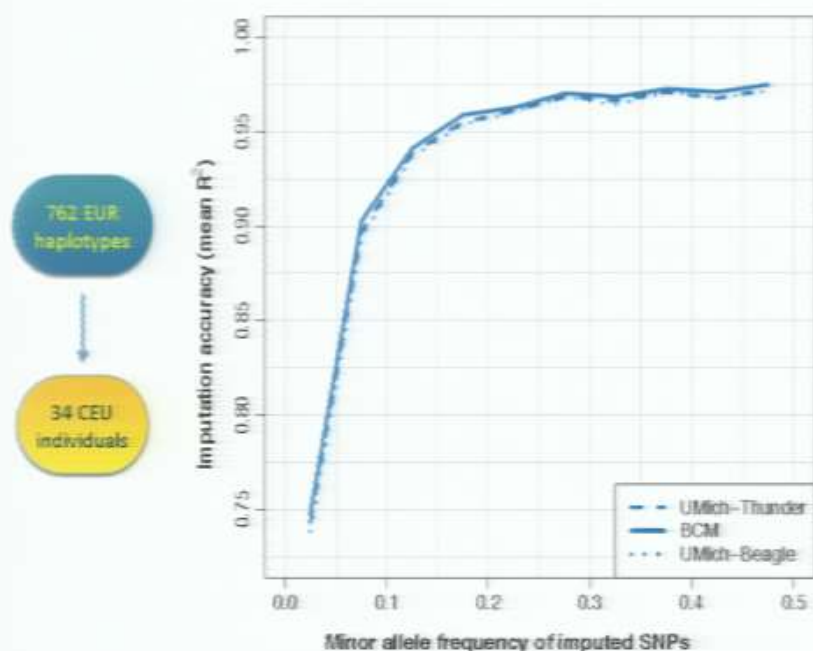
GWAS of imputed genotypes

- Increased power
- Better resolution
- Facilitates meta-analysis



1000 Genomes Project

Imputation Accuracy



Improvements from pilot to phase 1 using GAIN psoriasis dataset

Reference	MAF 1-3%	MAF 3-5%	MAF >5%
1000G Pilot	.69	.77	.91
1000G Phase I	.82	.85	.94

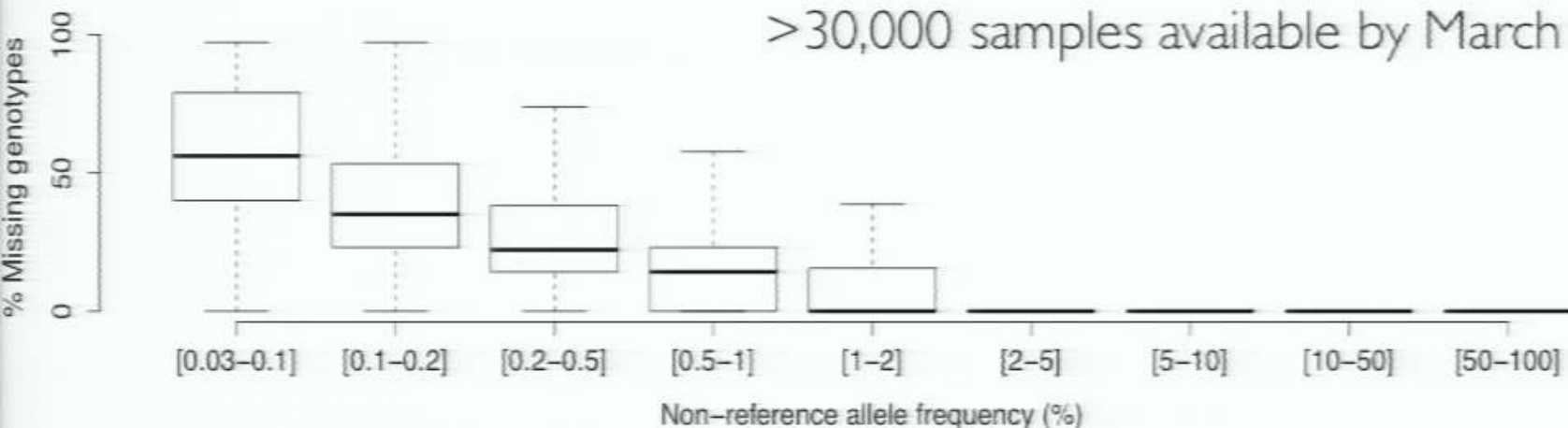
The Haplotype Consortium: Chromosome 20 pilot

COHORT	#samples	coverage	#SNPs
1000GP	379	4	338,766
UK10K	3,781	6.5	1,004,955
ORCADES	399	4	259,373
FINNS	1,941	4	315,539
GoNL	748	12	434,984
GoT2D	2,874	4	570,847
AMD	630	4	355,438
SARDINIA	2,120	4	372,632
MTFS	687	6	468,672
TOTALS	13,559		1,649,648

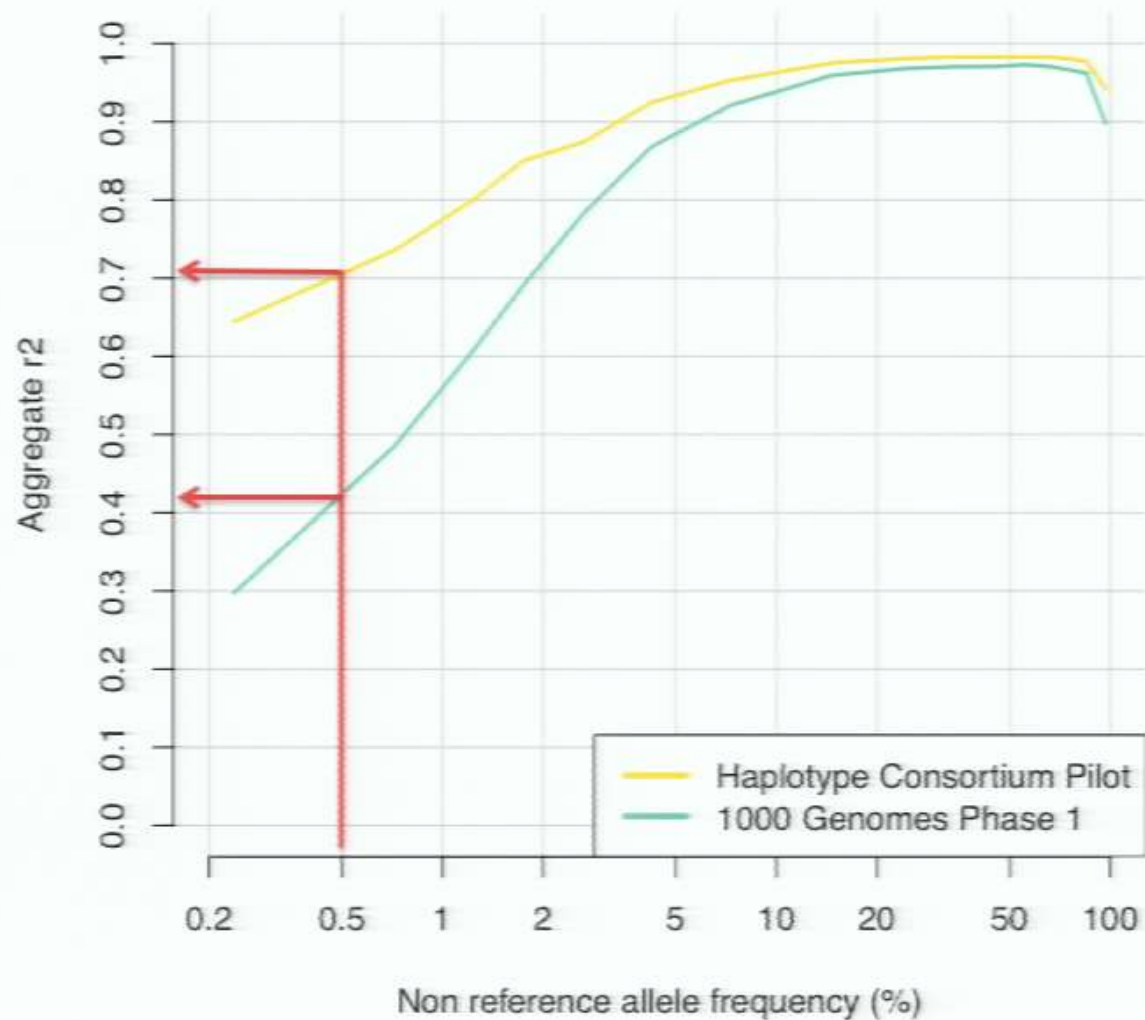
We created a union set of sites across 9 studies

Only use sites with > 10 minor alleles
i.e. $MAF > 0.03\%$

- 427,589 sites
- 26.5% missing genotypes:
- fill in by imputation



Pilot Results: 13,559 sequences chr20



Haplotype Consortium

Jonathan Marchini and Olivier Delaneau, Oxford

1000 Genomes

Goncalo Abecasis
Gil McVean

Sardinia

Francesco Cucca
Serena Sanna
David Schlessinger
Carlo Sidore

Finns

Aarno Palotie

GoNL

Cisca Wijmenga
Paul I.W. de Bakker
Morris A. Swertz
Androniki Menelaou

UK10K

Nicole Soranzo
Nick Timpson
George Davey-Smith
Tim Spector

Orcades

Jim Wilson

AMD

Anand Swaroop
Dwight Stambolian
Emily Chew
Xiaowei Zhan

HUNT

Cristen Willer
Kristian Hveem

GoT2D

Mark McCarthy
David Altshuler
Mike Boehnke

MTFS

Scott Vrieze
Matt McGue
Bill Iacono

Bipolar

Mike Boehnke
Richard Myers

Helic

Eleftheria Zeggini

Crohns/UC

Jeff Barrett
Carl Anderson

Italian Isolates

Paolo Gasparini
Nicole Soranzo
Daniela Toniolo
Nicola Pirastu

Oxford

Warren Kretzschmar

WT Sanger Institute

Shane McCarthy

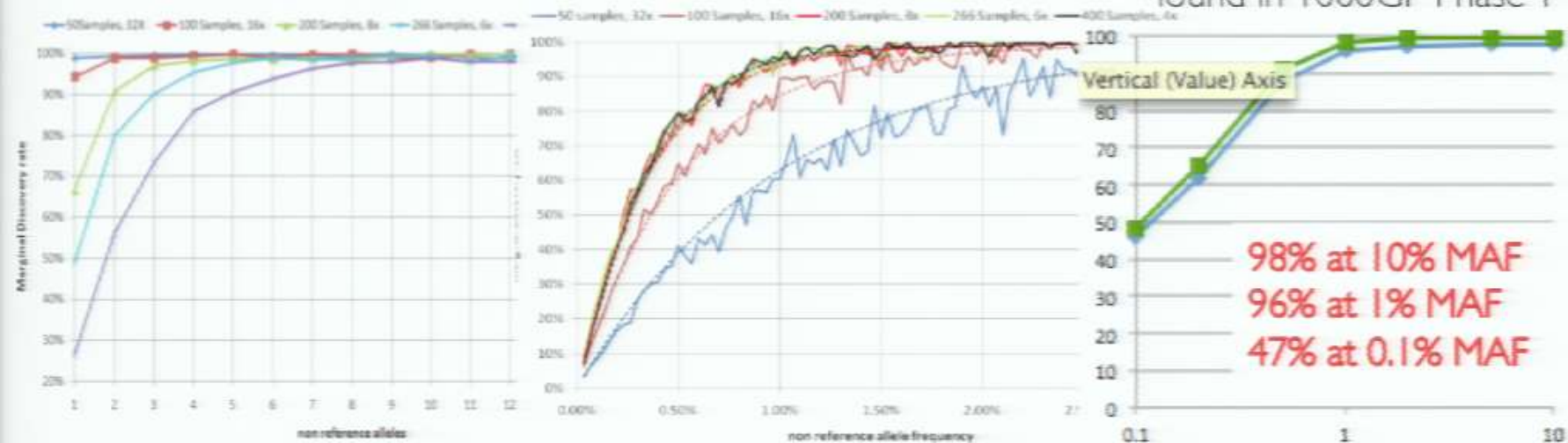
Michigan

Christian Fuchsberger
Hyun Min Kang

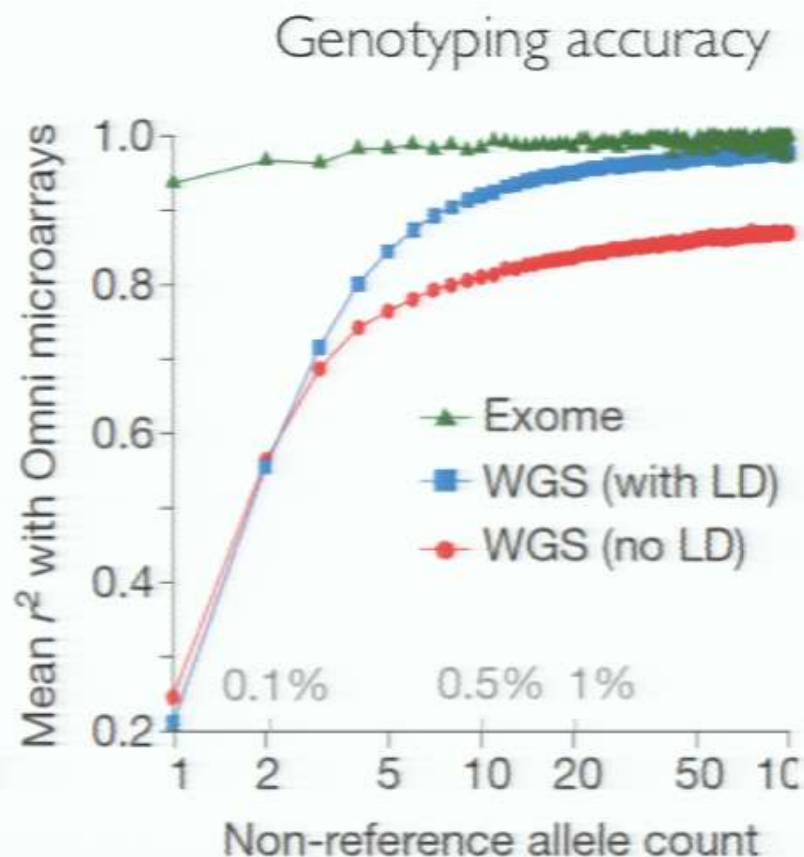
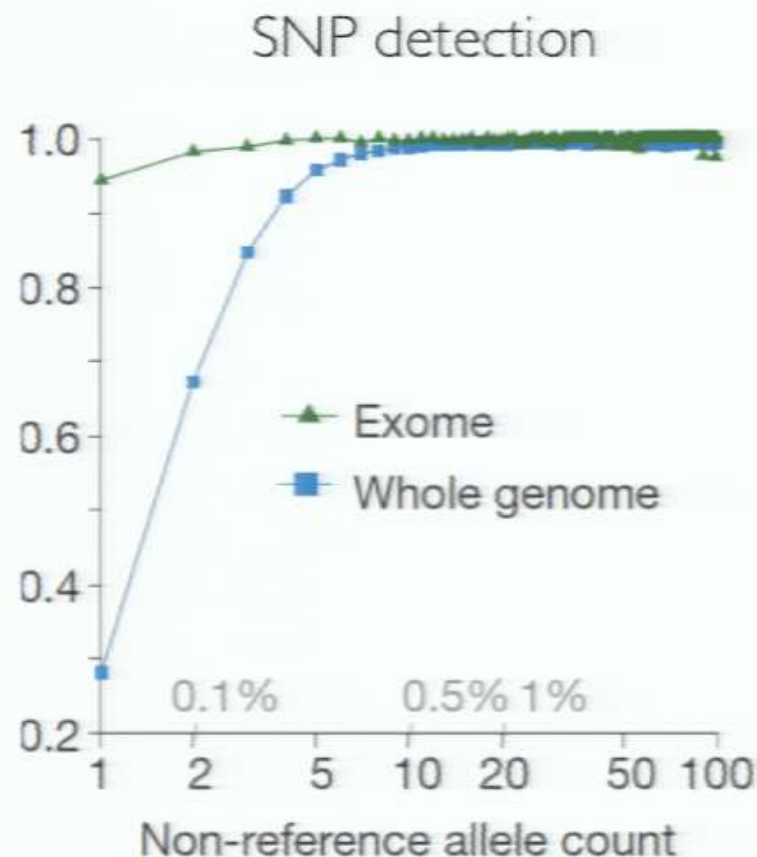
Imputation is also used to integrate data from low coverage sequencing

- To find the sequence of a single sample stand alone one needs to sequence at $\sim 30\times$ depth or more
- To find low frequency variants we want to sequence many samples
- Spread sequence across more samples

Fraction of SNPs discovered in 2453 UK10K samples found in 1000GP Phase I



Phase I power and genotyping accuracy



How to improve/speed up imputation?

- HMM based methods such as IMPUTE are computationally heavy
 - Quadratic (linear) in number of reference sequences
 - Typically use MCMC sampling
 - 1000s of CPU days for current data sets
- They will not scale to millions of genomes

Updating sort order is linear time

$$y^k[k] \quad y^{k+1}[k+1]$$
[illegible]

Reverse sorted prefixes at k

Let u_i be the number of 0s in y before i , i.e. $u_i = i - \sum_{j < i} y_j$

And c be the total number of 1s in y , i.e. $c = \sum y_j$

Then i maps to u_i if $y_i = 0$, and $c + i - u_i$ if $y_i = 1$.

New data structure for fast matching

Positional Burrows-Wheeler Transform (PBWT)

PBW[k]

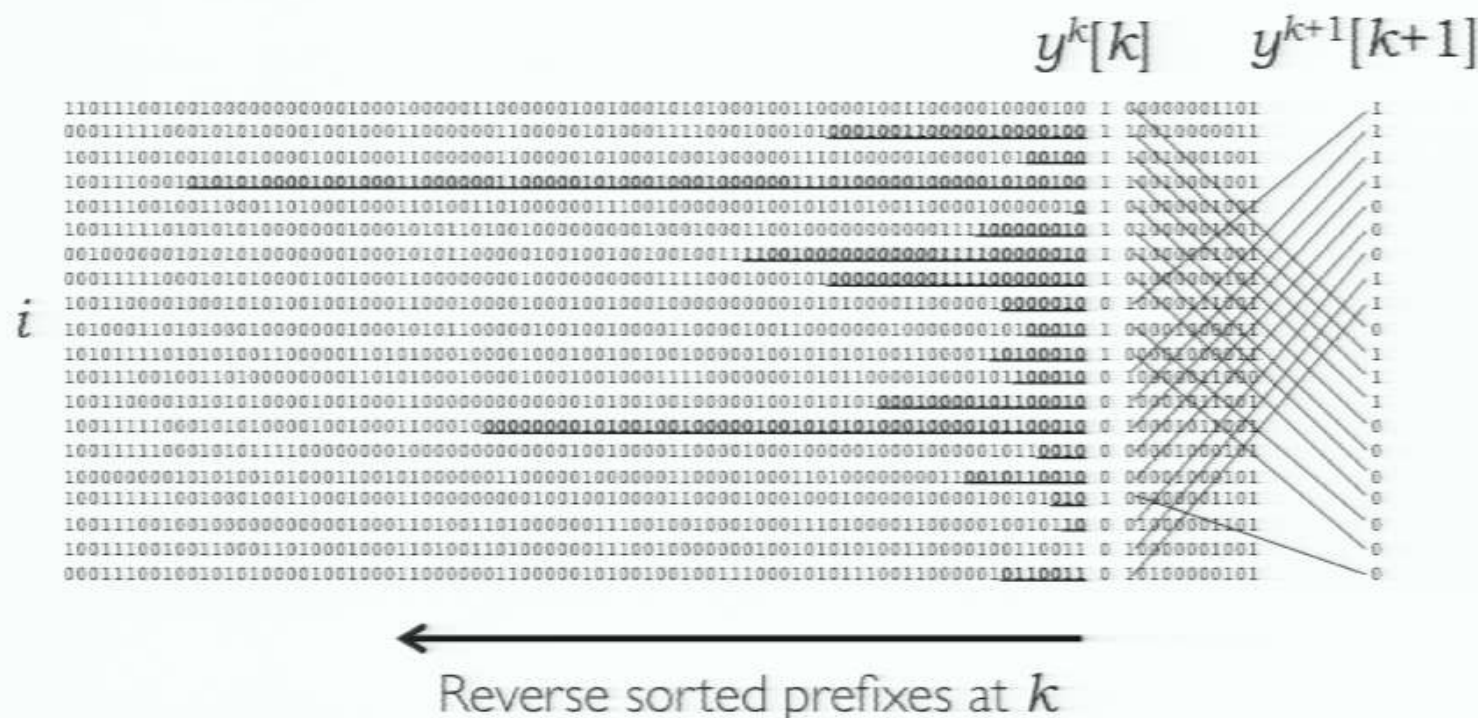
```
01001010100000011101110010010000000000100010000011000000100100010101000100110000100110000010000100 1 0000000110110010000
010010010110110110011111000101010000100100011000000110000010100011110001000101000100110000010000100 1 1001000001100011110
0110101100010000100111001001010100001001000110000001100000101000100010000001110100000100000101000100 1 1001000100100100010
0100100100010000100111000101010100001001001101001101000000111001000000100101010100110000100000010 1 1001000100100100010
01100101000100001001110010011000110100010001101001101000000111001000000100101010100110000100000010 1 0100000100100010010
01000101010100001001111101010101000000010001010110100100000000010001000110010000000000111100000010 1 0100000100100010000
01100101000100110010000001010101000000010001010110000010010010010010011110010000000000011100000010 1 0100000100100010010
1100010100010011000111110001010100001001000110000000010000000000111100010001010000000011100000010 1 0100000010101000010
010010110001000010011000010001010100100100011000100001000100100010000000000101010000110000010000010 0 1000011100100010000
011001010111000010100011010100010000000100010101100000100100100001100001001100000001000000010100010 1 0000100001100010000
0110010101100000101011110101000110000011010100010000100010010010010000010010101010001100001000010 1 0000100001100010010
01100101000100001001110010011000001010101000010010001100000000000010100100100000100101010001000010 0 100001100000100010
010010010000110010011000010101010000100100011000000000000010100100100000100101010100010000101100010 0 1000101100100000010
0110010100010101100111110001010100001001000110001000000001000100100000100101010100010000100100010 0 1000101100100000010
0100100101101101100111110001010111000000001000000000000100100001100001000100000100010000010110010 0 0000100010110000110
010010010110000010000000010101001010001100101000000110000010000001100001000110100000000110010110010 0 0000100010110000111
0000010100100011100111110010001001100010001100000000010010010000110000100010000100001000100101010 1 0000000110100010000
0100100101101101100111001001000000000010001101001101000000111001001000100011010000110000010010110 0 0100000110110010010
0100010100010000100111001001100011010001000110100110100000011100100000001001010100110000100110011 0 1000000100100000010
010010011010001100011100100101010000100100011000000110000010100100100111000101011100110000010110011 0 1010000010101000000
```

←
Reverse sorted prefixes at k

Matches are adjacent in the sort order

Analogous to BWT used by BWA, BowTie etc. for sequence matching

Updating sort order is linear time

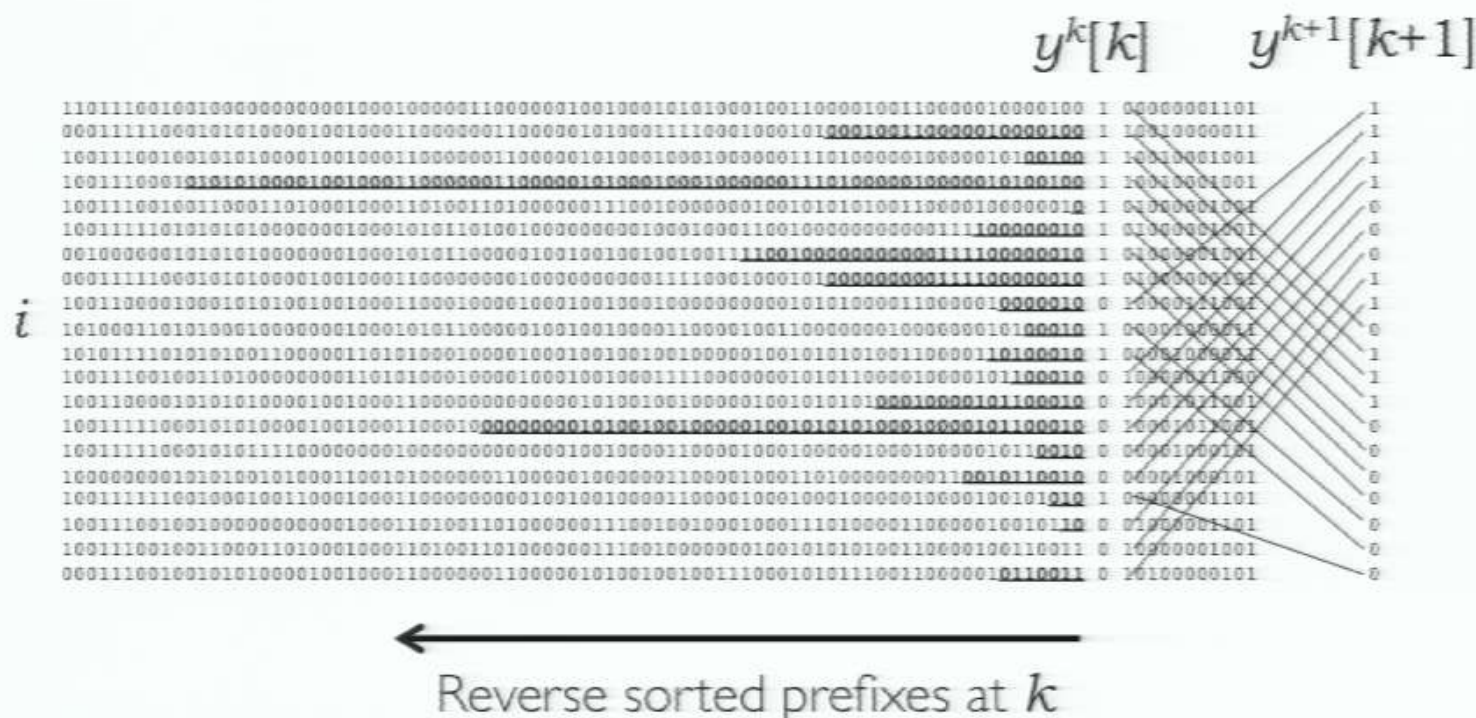


Let u_i be the number of 0s in y before i , i.e. $u_i = i - \sum_{j < i} y_j$

And c be the total number of 1s in y , i.e. $c = \sum y_j$

Then i maps to u_i if $y_i = 0$, and $c + i - u_i$ if $y_i = 1$.

Updating sort order is linear time



Let u_i be the number of 0s in y before i , i.e. $u_i = i - \sum_{j < i} y_j$

And c be the total number of 1s in y , i.e. $c = \sum y_j$

Then i maps to u_i if $y_i = 0$, and $c + i - u_i$ if $y_i = 1$.

Data compression – run length encoding

- Encode runs in bytes using
 - 1 bit for the value
 - 2 bits select run length units: 1 or 64 or 2048
 - 5(6) bits give number of units
- Simulate 100k sequences of length 20Mbps with ARG simulator MaCS (Chen et al. 2009)
 - 370,264 sites (one per 54bp): 37GB raw output
 - gzip compresses to 1.02GB (~35x compression)
 - PBWT compresses to 7.7MB (~4800x compression)
 - Native order run length encodes to ~2GB

Updating sort order is linear time

$$y^k[k] \quad y^{k+1}[k+1]$$
Reverse sorted prefixes at k

Let u_i be the number of 0s in y before i , i.e. $u_i = i - \sum_{j < i} y_j$

And c be the total number of 1s in y , i.e. $c = \sum y_j$

Then i maps to u_i if $y_i = 0$, and $c + i - u_i$ if $y_i = 1$.

Data compression – run length encoding

- Encode runs in bytes using
 - 1 bit for the value
 - 2 bits select run length units: 1 or 64 or 2048
 - 5(6) bits give number of units
- Simulate 100k sequences of length 20Mbps with ARG simulator MaCS (Chen et al. 2009)
 - 370,264 sites (one per 54bp): 37GB raw output
 - gzip compresses to 1.02GB (~35x compression)
 - PBWT compresses to 7.7MB (~4800x compresssion)
 - Native order run length encodes to ~2GB

Compression performance

Simulated data subsets

	1000	2000	5000	10000	20000	50000	100000
PBWT	1685629	1956360	2783732	3372188	4145516	5688290	7697194
haps.gz	10515008	21340464	53246970	105558782	209332249	517432833	1024614062
factor	6.2	10.9	19.1	31.3	50.5	91.0	133.1
bytes/site	4.6	5.3	7.5	9.1	11.2	15.4	20.8

“Real” data: 1000 Genomes phase I chromosome I

2184 chromosomes, 3007196 sites

PBWT **51186641**

gzip **302883517**

factor **5.9**

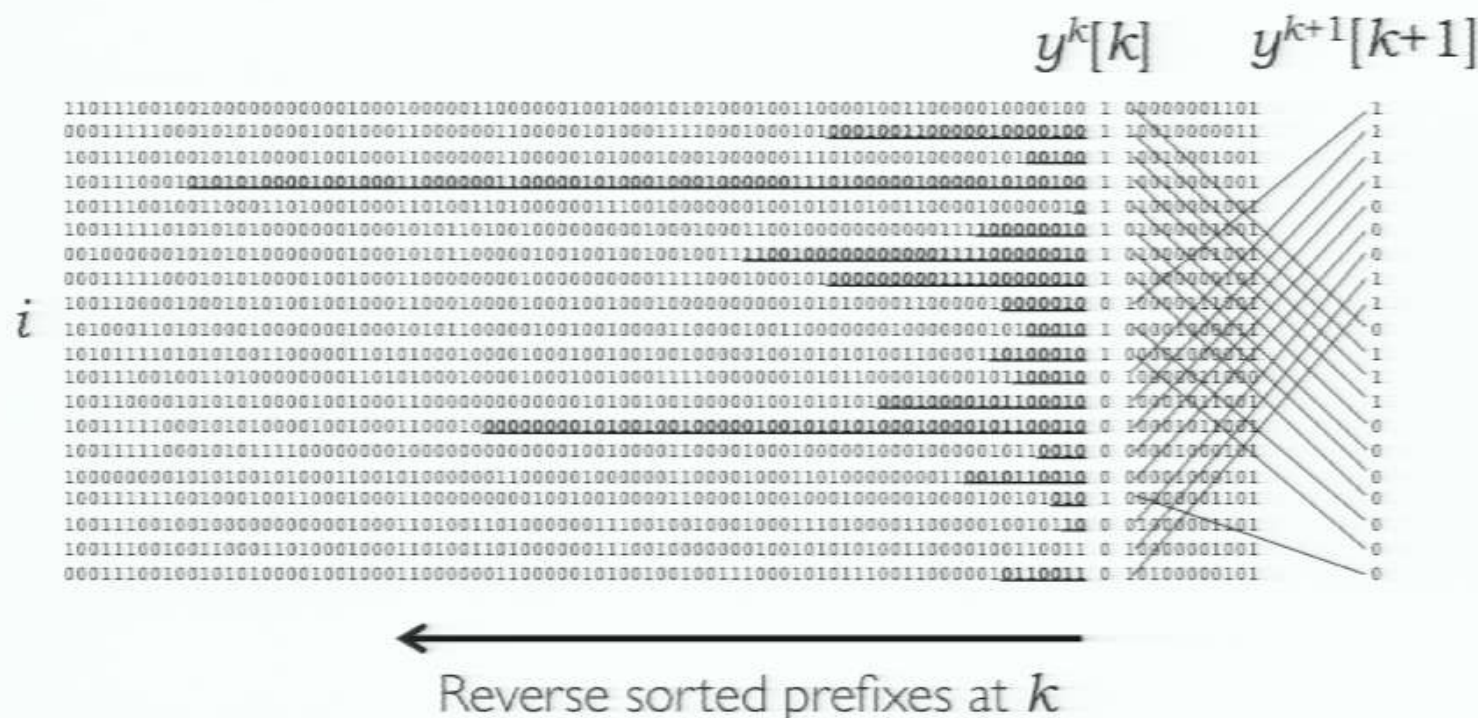
Maximal matches of new sequences to a reference panel

- Update rule the same as for sort update
- Option 1: build indexes on top of PBWT
 - Very fast, effectively independent of the size of the reference panel so $O(N)$ time
 - Quite memory hungry, $\sim 13NM$ bytes
- Option 2: match a batch of new sequences as you pass once through the PBWT
 - $O(NM)$ time, but very memory efficient

Data compression – run length encoding

- Encode runs in bytes using
 - 1 bit for the value
 - 2 bits select run length units: 1 or 64 or 2048
 - 5(6) bits give number of units
- Simulate 100k sequences of length 20Mbps with ARG simulator MaCS (Chen et al. 2009)
 - 370,264 sites (one per 54bp): 37GB raw output
 - gzip compresses to 1.02GB (~35x compression)
 - PBWT compresses to 7.7MB (~4800x compresssion)
 - Native order run length encodes to ~2GB

Updating sort order is linear time



Let u_i be the number of 0s in y before i , i.e. $u_i = i - \sum_{j < i} y_j$

And c be the total number of 1s in y , i.e. $c = \sum y_j$

Then i maps to u_i if $y_i = 0$, and $c + i - u_i$ if $y_i = 1$.

Data compression – run length encoding

- Encode runs in bytes using
 - 1 bit for the value
 - 2 bits select run length units: 1 or 64 or 2048
 - 5(6) bits give number of units
- Simulate 100k sequences of length 20Mbp with ARG simulator MaCS (Chen et al. 2009)
 - 370,264 sites (one per 54bp): 37GB raw output
 - gzip compresses to 1.02GB (~35x compression)
 - PBWT compresses to 7.7MB (~4800x compresssion)
 - Native order run length encodes to ~2GB

Compression performance

Simulated data subsets

	1000	2000	5000	10000	20000	50000	100000
PBWT	1685629	1956360	2783732	3372188	4145516	5688290	7697194
haps.gz	10515008	21340464	53246970	105558782	209332249	517432833	1024614062
factor	6.2	10.9	19.1	31.3	50.5	91.0	133.1
bytes/site	4.6	5.3	7.5	9.1	11.2	15.4	20.8

“Real” data: 1000 Genomes phase I chromosome I

2184 chromosomes, 3007196 sites

PBWT **51186641**

gzip **302883517**

factor **5.9**

Maximal matches of new sequences to a reference panel

- Update rule the same as for sort update
- Option 1: build indexes on top of PBWT
 - Very fast, effectively independent of the size of the reference panel so $O(N)$ time
 - Quite memory hungry, $\sim 13NM$ bytes
- Option 2: match a batch of new sequences as you pass once through the PBWT
 - $O(NM)$ time, but very memory efficient

Time to match 1000 new sequences

Reference panel size	1,000	2,000	5,000	10,000	20,000	50,000
Naïve	52.1	103.8	258.9	519.2	1035.5	2582.6
Indexed	0.9u	0.9u	0.9u	0.9u	1.1u	1.7u
	0.1s	0.1s	0.1s	0.2s	0.5s	15s
Batch	2.3u	2.5u	3.5u	4.8u	6.8u	12.1u
	0.1s	0.1s	0.1s	0.1s	0.1s	0.1s

The reference panel here is pseudo-genotype array data, made of 10% of sites with $MAF > 0.05$ from the full sequence simulation.

The naïve method compares each sequence to each previous sequence, inefficiently.

What about inference?

- Prediction is closely related to compression

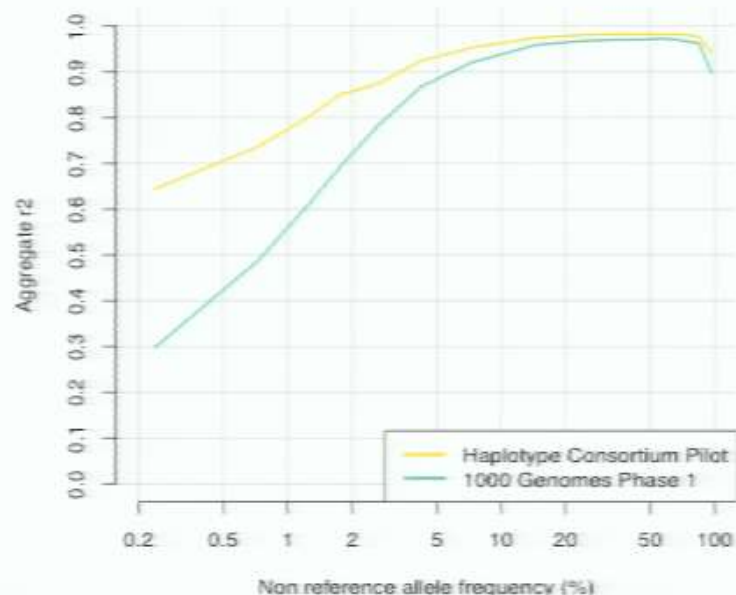
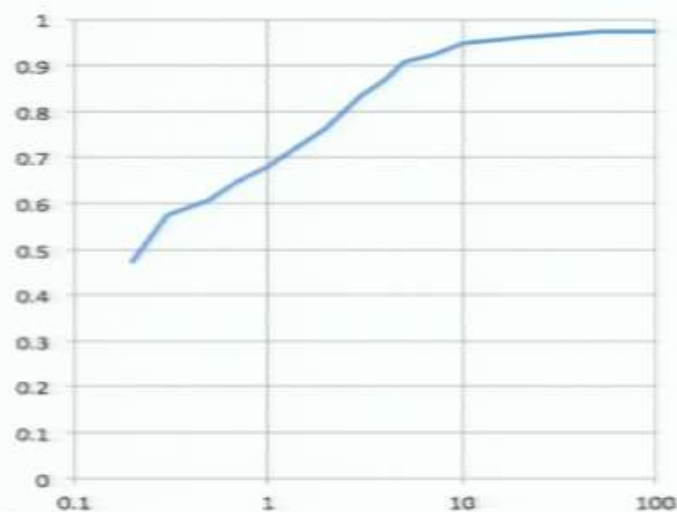
```
01001010100000011101110010010000000000100010000011000000100100010101000100110000100110000010000100 1
01001001011011011001111100010101000010010001100000011000001010001111000100010100010011000001000100 1
0110101100010000100111001001010100001001000110000001100000101000100010000001110100000100000101000100 1
0100100100010000100111000101010100001001000110000001100000101000100010000100001000010000100100100 1
0110010100010000100111001001100011010011010011010000001110010000000100101010100110000100000010 1
010001010101000010011111010101010000000100010101101001000000000100010001100100000000000111100000010 1
011001010001001100100000001010101000000010001010110000010010010010010011110010000000000111100000010 1
110001010001001100011111000101010000100100011000000001000000000011110001000101000000000111100000010 1
010010110001000010011000010001010100100100011000100001000100010000000000101010000110000010000010000010 0
011001011110000101000110101000100000001000101011000001001001000011000010011000000010000000101000010 1
01100101011000001010111101010100011000001101010001000010001001001000000100101010100110000110100010 1
011001010001000010011100100110100000000110101000100001000100100011110000000101011000010000101100010 0
010010010000110010011000010101010000100100011000000000000010100100100000100101010100010000101100010 0
011001010001010110011111000101010000100100011000100000000100100100000100101010100010000101100010 0
01001001011011011001111100010101111000000001000000000000100100001100001000100000100010000010110010 0
01001001011000001000000000101001010001100101000000110000010000001100001000110100000000110010110010 0
00000101001000111001111110010001001100001000110000000001001001000011000010001000010000010000100101010 1
010010010110110110011100100100000000000100011010011010000001110010010001000111010000110000010010110 0
010001010001000010011100100110001101000100011010011010000001110010000000100101010100110000100110011 0
010010011010001100011100100101010000100100011000000101001001001110001010111001100000101100110000010 0
```

y

- Consider a generative model for y given d
- e.g. Markov given by $p(y_{i+1} | y_i, d_i)$
- Iterate to generate whole sequence sets

Imputation (very preliminary)

- Empirically $p(y_{i+1} \neq y_i | d_i) \approx e^{-\alpha - \beta d}$
 - Expected given $d \sim$ Gumbel extreme value
- Imputation under Markov model is local, so very fast



26,618 haps, 428,131 sites (~60Mb)

25 mins on Mac air – results a bit worse than using IMPUTE

What about inference?

- Prediction is closely related to compression

```
01001010100000011101110010010000000000100010000011000000100100010101000100110000100110000010000100 1
0100100101101101100111110001010100001001000110000001100000101000111100010001010001001100000100 1
011010110001000010011100100101010000100100011000000110000010100010001000000111010000010000010100 1
0100100100010000100111000101010100001001000110000001100000101000100010000100001000010100100 1
0110010100010000100111001001100011010001000110100110100000011100100000010010101010011000010000001 1
01000101010100001001111101010101000000010001010110100100000000010001000110010000000000011110000001 1
011001010001001100100000010101010000000100010101100000100100100100100111100100000000001110000001 1
1100010100010011000111110001010100001001000110000000010000000000111100010001010000000011110000001 1
0100101100010000100110000100010101001001000110001000010001001000100000000010101000011000001000001 0
01100101011100001010001101010001000000010001010110000010010010000110000100110000000100000001010001 1
011001010110000010101111010101000110000011010100010000100010010010000010010101010011000011010001 1
01100101000100001001110010011010000000011010100010000100010010001111000000010101100001000010110001 0
01001001000011001001100001010101000010010001100000000000001010010010000010010101010001000010110001 0
01100101000101011001111100010101000010010001100010000000010010010000010010101010001000010110001 0
010010010110110110011111000101011100000000100000000000010010000110000100010000010001000001011001 0
01001001011000001000000001010100101000110010100000011000001000000110000100011010000000011001011001 0
00000101001000111001111110010001001100010001100000000010010010000110000100010001000001000010010101 1
010010010110110110011100100100000000000100011010011010000001100100100010001110100001100000100101 0
010001010001000010011100100110001101000100011010011010000001100100000001001010100110000100110011 0
01001001101000110001110010010101000010010001100000010100100100111000101011100110000010110011 0
```

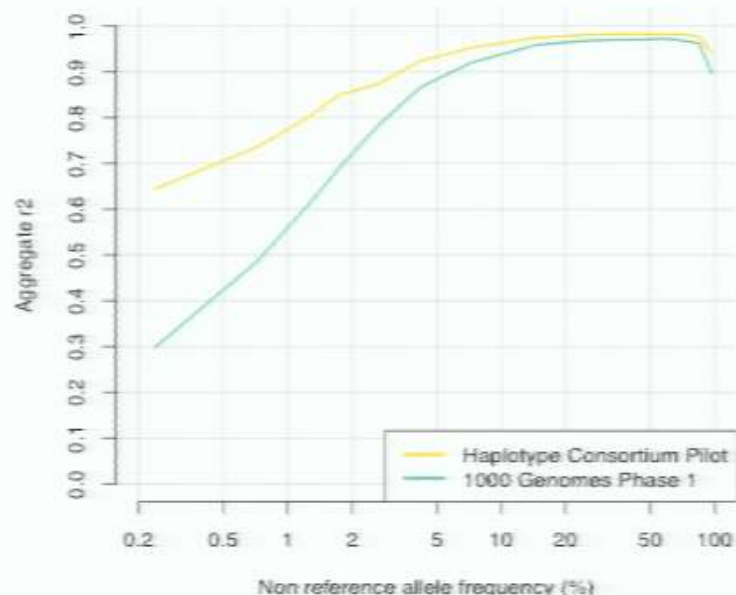
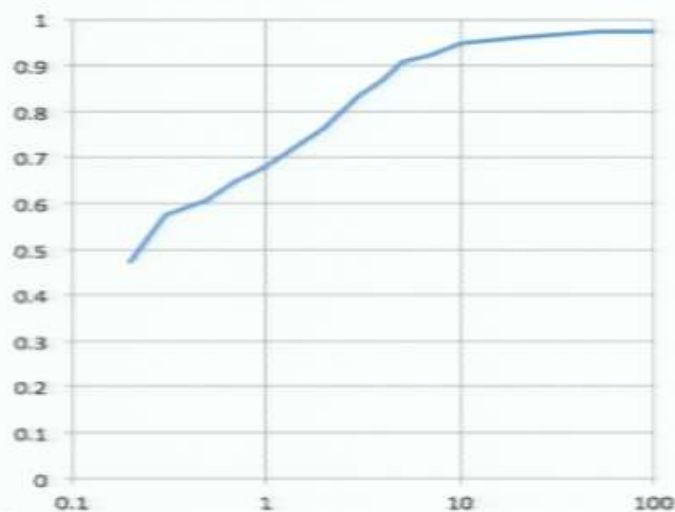
y

d

- Consider a generative model for y given d
- e.g. Markov given by $p(y_{i+1} | y_i, d_i)$
- Iterate to generate whole sequence sets

Imputation (very preliminary)

- Empirically $p(y_{i+1} \neq y_i | d_i) \approx e^{-\alpha - \beta d}$
 - Expected given $d \sim$ Gumbel extreme value
- Imputation under Markov model is local, so very fast



26,618 haps, 428,131 sites (~60Mb)

25 mins on Mac air – results a bit worse than using IMPUTE

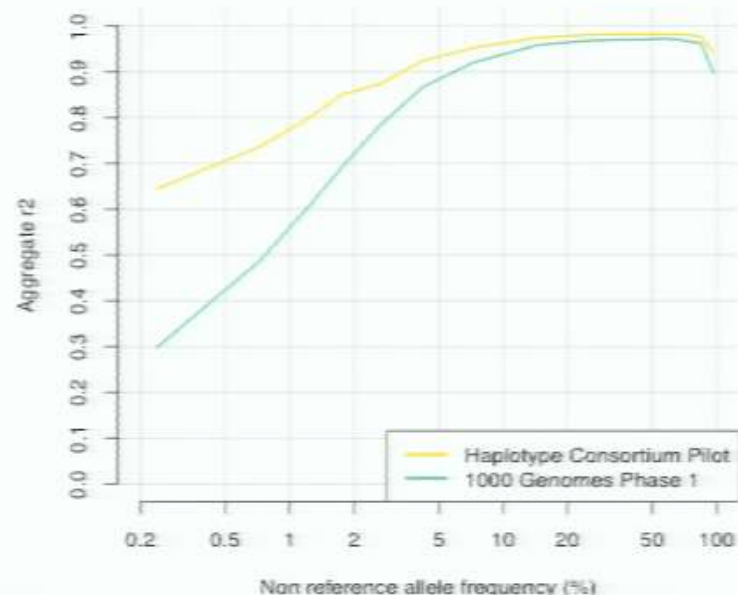
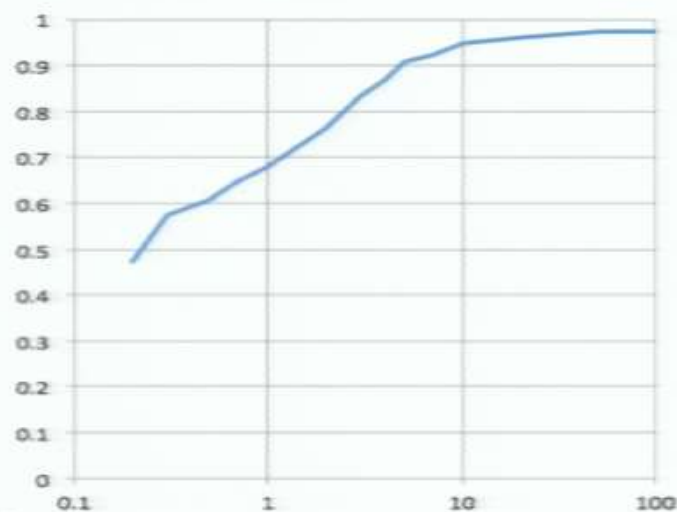
Alternative model: tree defined by d_i ?



- Potentially attractive relationship to ARG
- But data not consistent with tree in hard form
 - Sample trees consistent with data weighting via d_i

Imputation (very preliminary)

- Empirically $p(y_{i+1} \neq y_i | d_i) \approx e^{-\alpha - \beta d}$
 - Expected given $d \sim$ Gumbel extreme value
- Imputation under Markov model is local, so very fast



26,618 haps, 428,131 sites (~60Mb)

25 mins on Mac air – results a bit worse than using IMPUTE

Back to the raw data

BW methods help there too

- Collect shotgun sequencing reads
 - Random fragments from the whole genome
 - Can be enriched e.g. for exome
- Map the reads to the reference genome
 - Potential problems in repetitive areas
 - Potential alignment problems
- Detect variants based on the multiple alignment of reads
 - Statistical issues allowing for errors and sampling

Many modern mappers (e.g. bwa) use the standard Burrows-Wheeler Transform

Non-unique match is an *interval* in Suffix Array

mississippi

i SA[i]

0	11	\$
1	10	i\$
2	7	ippi\$
3	4	issippi\$
4	1	ississippi\$
5	0	mississippi\$
6	9	pi\$
7	8	ppi\$
8	6	sippi\$
9	3	sissippi\$
10	5	ssippi\$
11	2	ssissippi\$

$L(si)=8$, $U(si)=10$



Many modern mappers (e.g. bwa) use the standard Burrows-Wheeler Transform

Non-unique match is an *interval* in Suffix Array

mississippi

i SA[i]

0	11	\$	←
1	10	i\$	←
2	7	ippi\$	
3	4	issippi\$	
4	1	ississippi\$	
5	0	mississippi\$	←
6	9	pi\$	
7	8	ppi\$	
8	6	sippi\$	←
9	3	sissippi\$	
10	5	ssippi\$	←
11	2	ssissippi\$	← i si

$L(si)=8, U(si)=10$

FM approach to matching:

update $L()$ and $U()$ as you extend the search string

$L()=0, U()=12$

$L(i)=1, U(i)=5$

$L(si)=8, U(si)=10$

Constant time update gives $O(M)$ search

Update uses FM-index

Ferragina and Manzini 2000

mississippi

i SA[i]

0	11	i \$
1	10	p i \$
2	7	s i p p i \$
3	4	s i s s i p p i \$
4	1	m i s s i s s i p p i \$
5	0	\$ m i s s i s s i p p i \$
6	9	p p i \$
7	8	i p p i \$
8	6	s s i p p i \$
9	3	s s i s s i p p i \$
10	5	i s s i p p i \$
11	2	i s s i s s i p p i \$

BW[i]

$$L(aS) = C[a] + P[a, L(S)]$$

$$U(aS) = C[a] + P[a, U(S)]$$

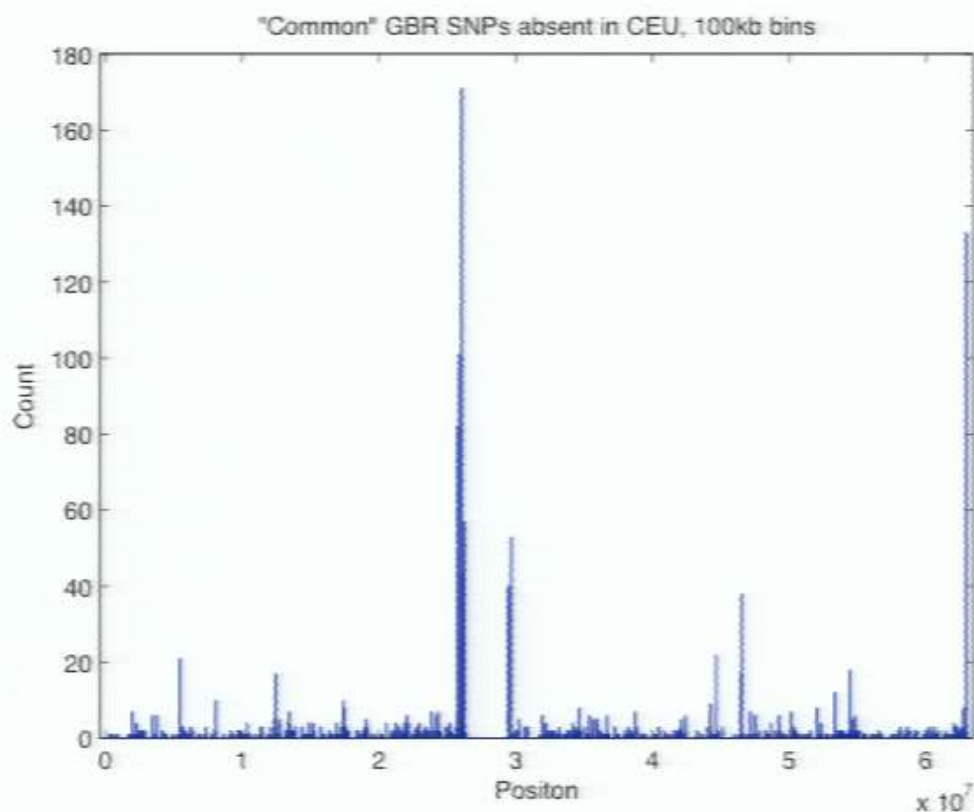
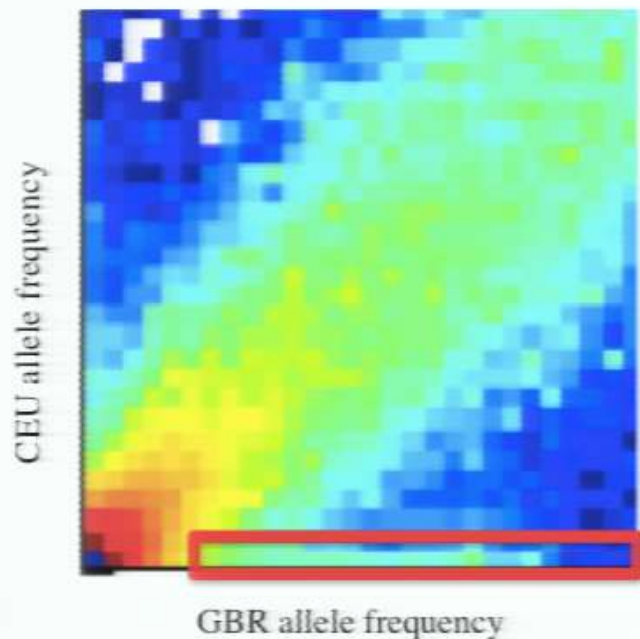
Where

$C[a]$ is the number of letters less than a in the target string X

$P[a, i]$ is the number of times a occurs in $BW[j]$ for $j < i$

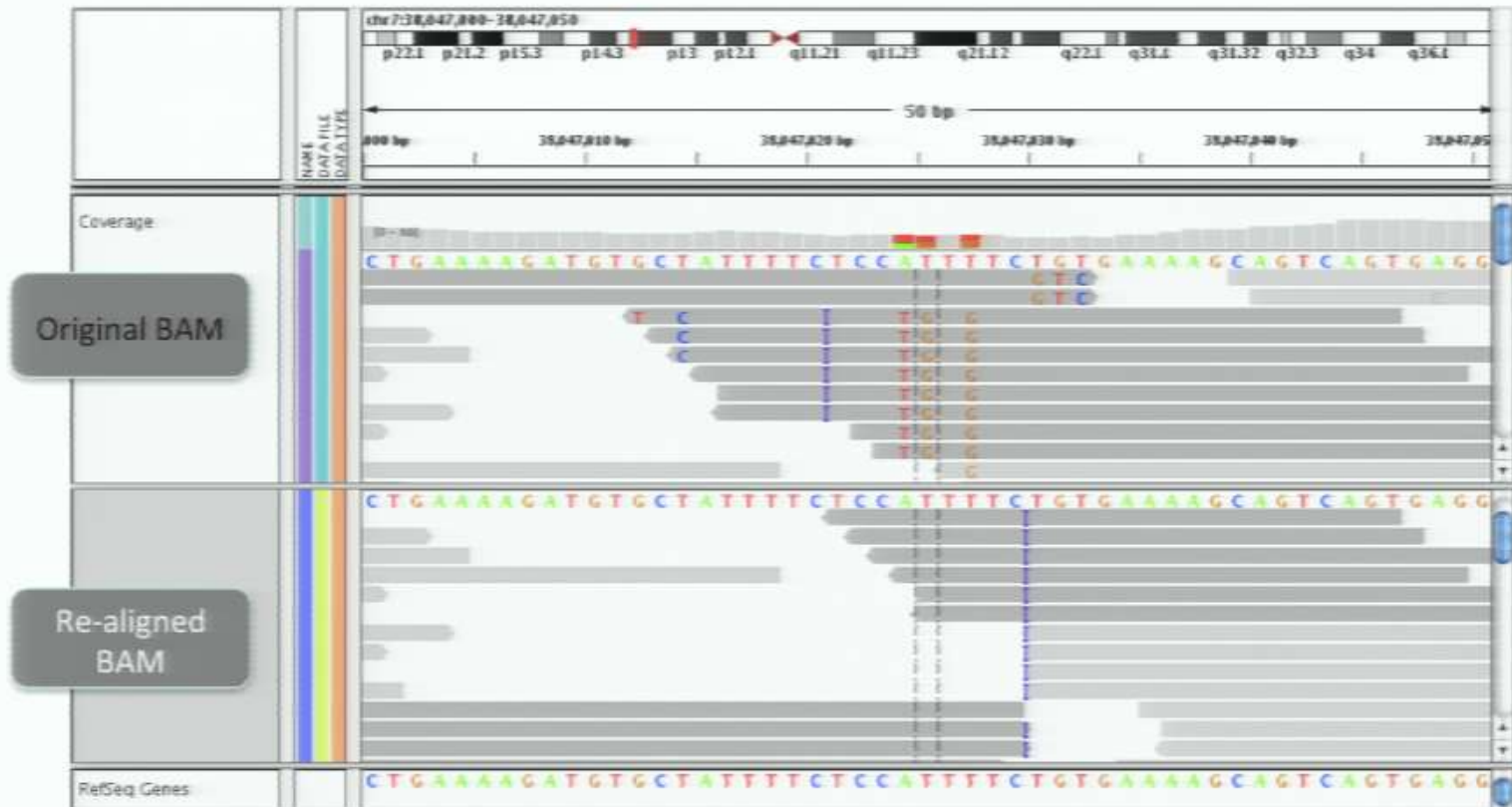
$BW[i] = X[SA[i]-1]$ is the Burrows-Wheeler transform

Errors due to mapping problems



Rare, population specific "SNPs" are clustered near centromere
Many of these are likely to be artefacts, e.g. ~1% of the total

Errors due to alignment at indels



- Control using local alignment uncertainty (BAQ: Li Heng), realignment, reassembly ...

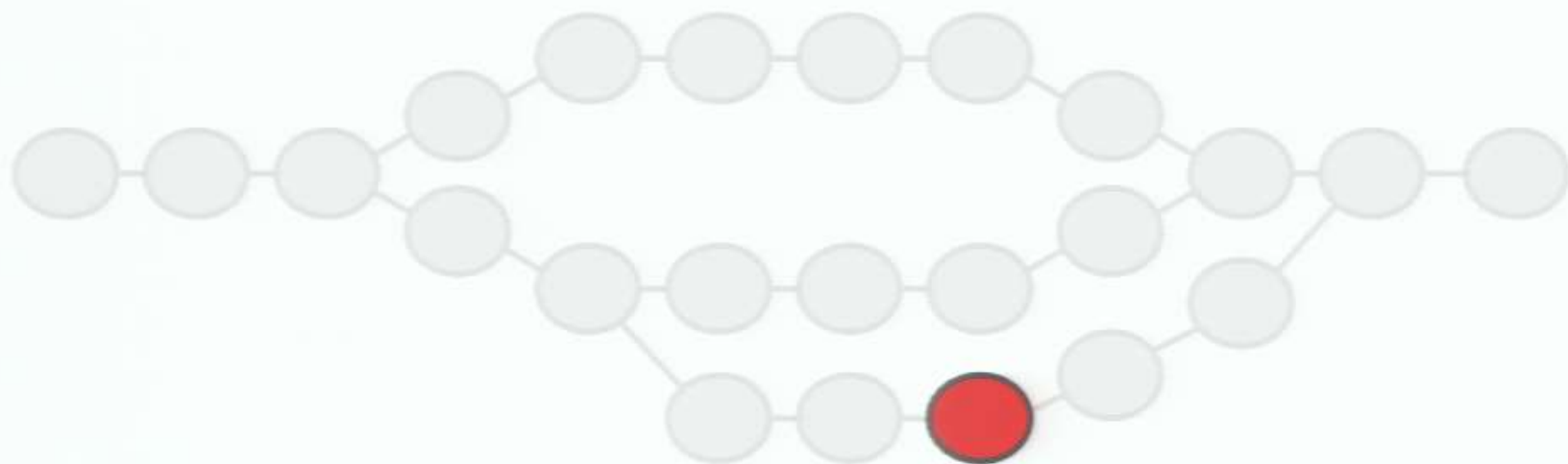
We are mapping to the wrong reference!

Genetic variation by assembly

- Reference free variation detection
 - De novo mutations by comparing child to parents
 - Somatic cancer mutations by comparing tumour to normal
 - Population variation by identifying segregating sequence
 - Individual variation by comparing to reference
- Could assemble first then compare contig sets
 - Time consuming
- We only want to assemble the *differences* between the samples

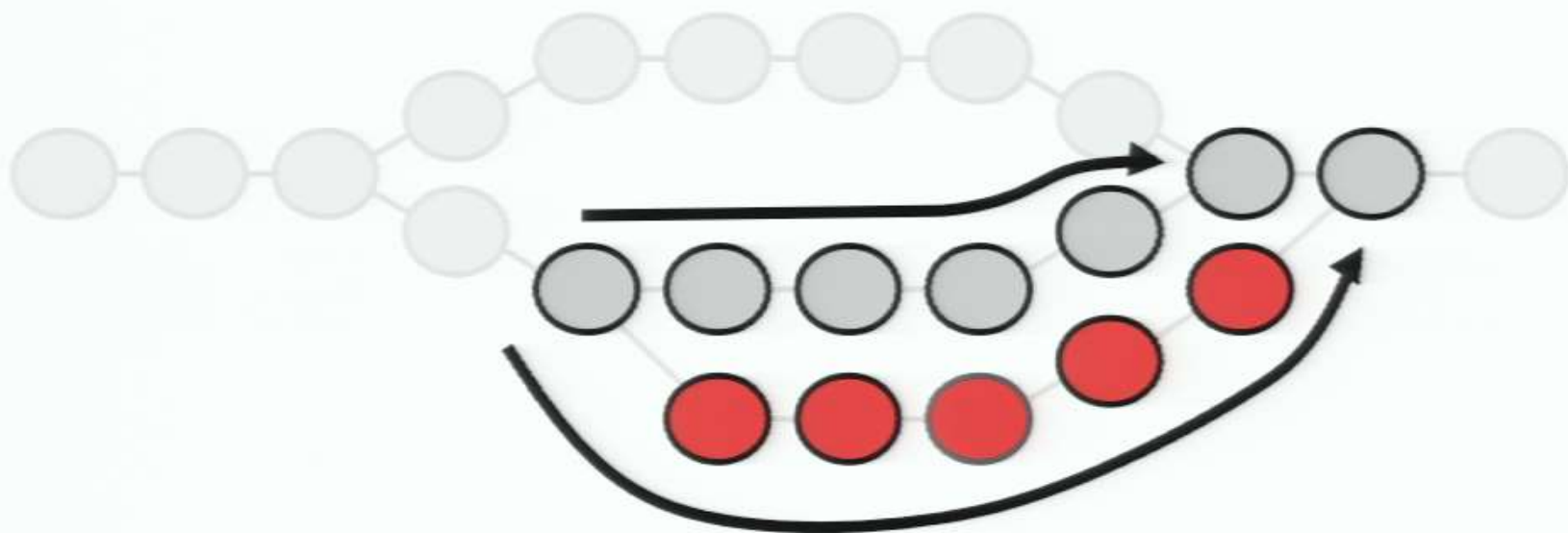
Genetic variation by assembly

Find unique reads/k-mers, then *locally* construct the string graph around reads containing these k-mers



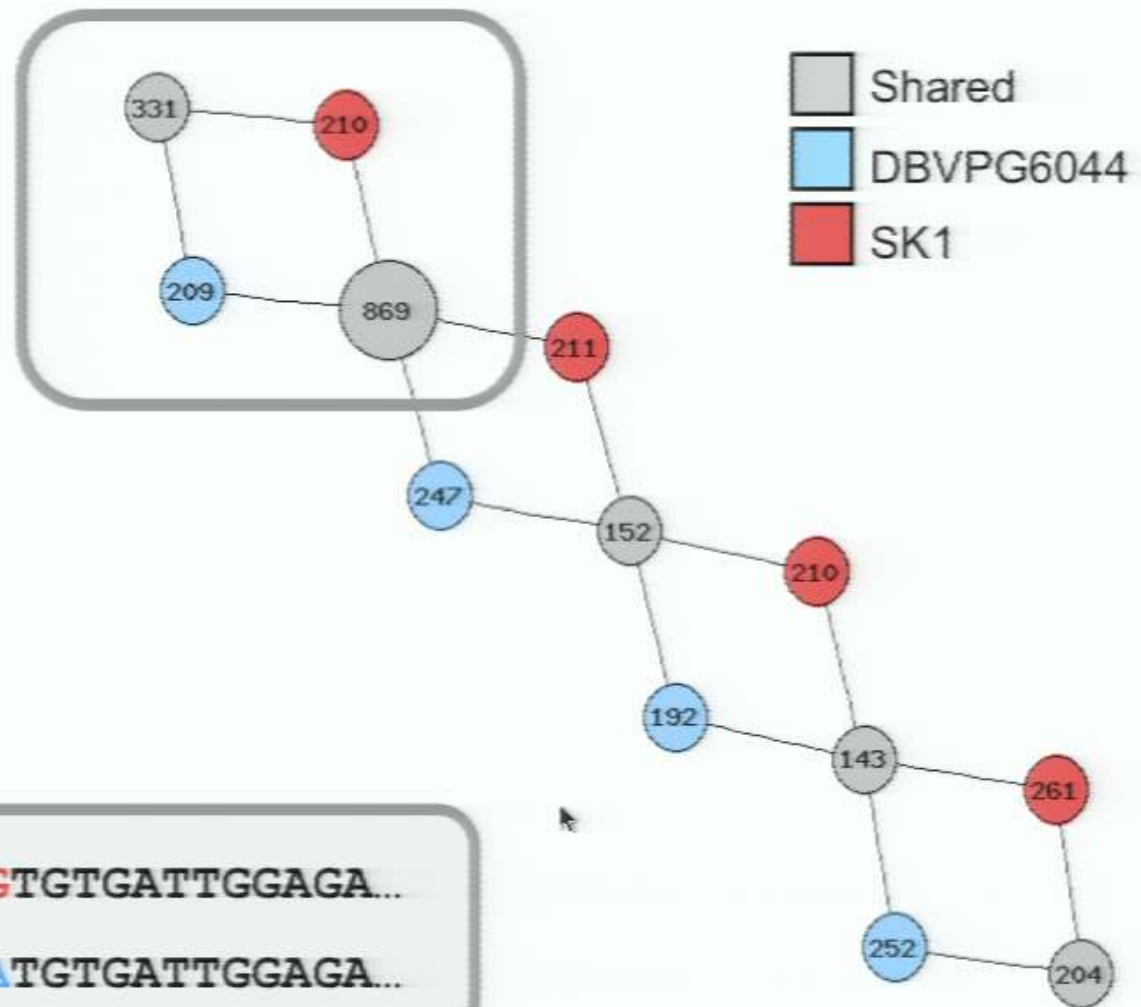
Genetic variation by assembly

...GATACATACGATCGATCTATGCATCCGAGCTGCTAGCTATCGATT... parent
...GATACATACGATCGATCTATGCATCCGAGC**A**GCTAGCTATCGATT... child



Evaluate evidence for haplotypes with Dindel (Albers et al. 2010) Bayesian methods

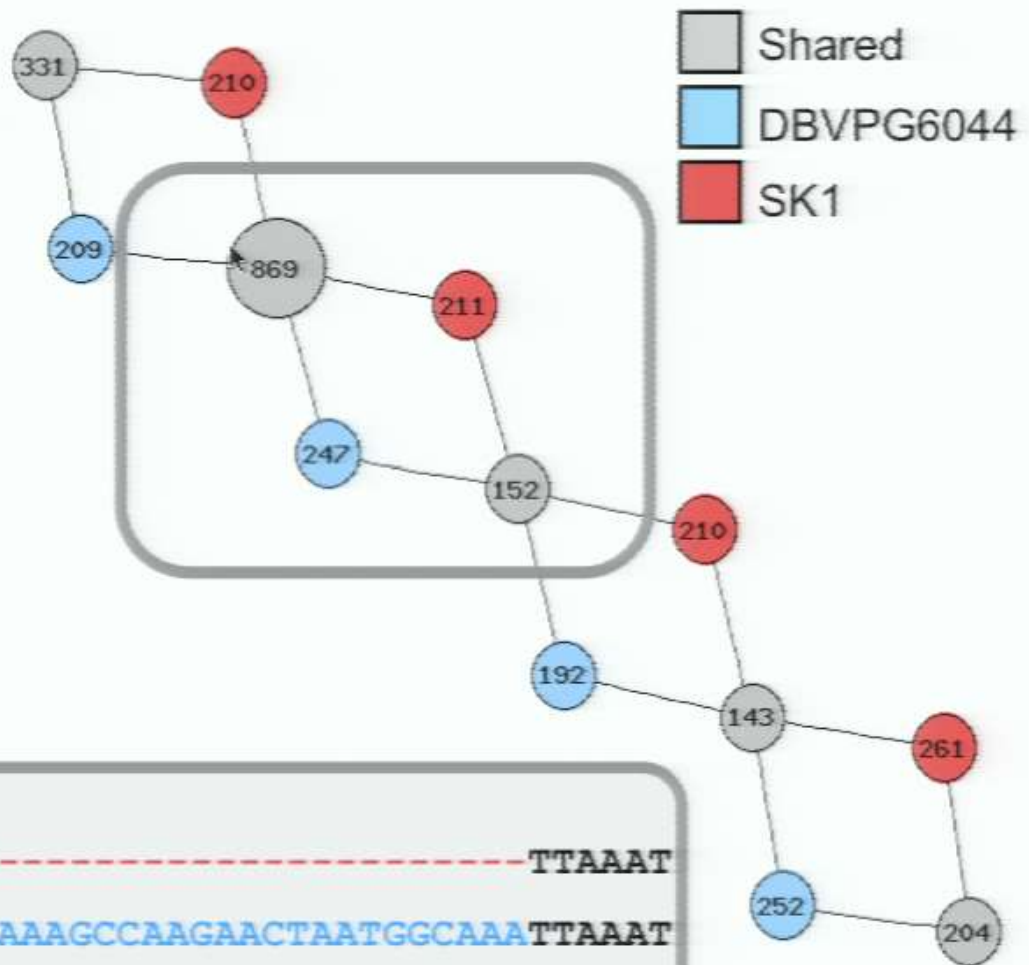
Two yeast strains co-assembled



...AAAGATATGCGC**GT**GTGATTGGAGA...

...AAAGATATGCGC**AT**GTGATTGGAGA...

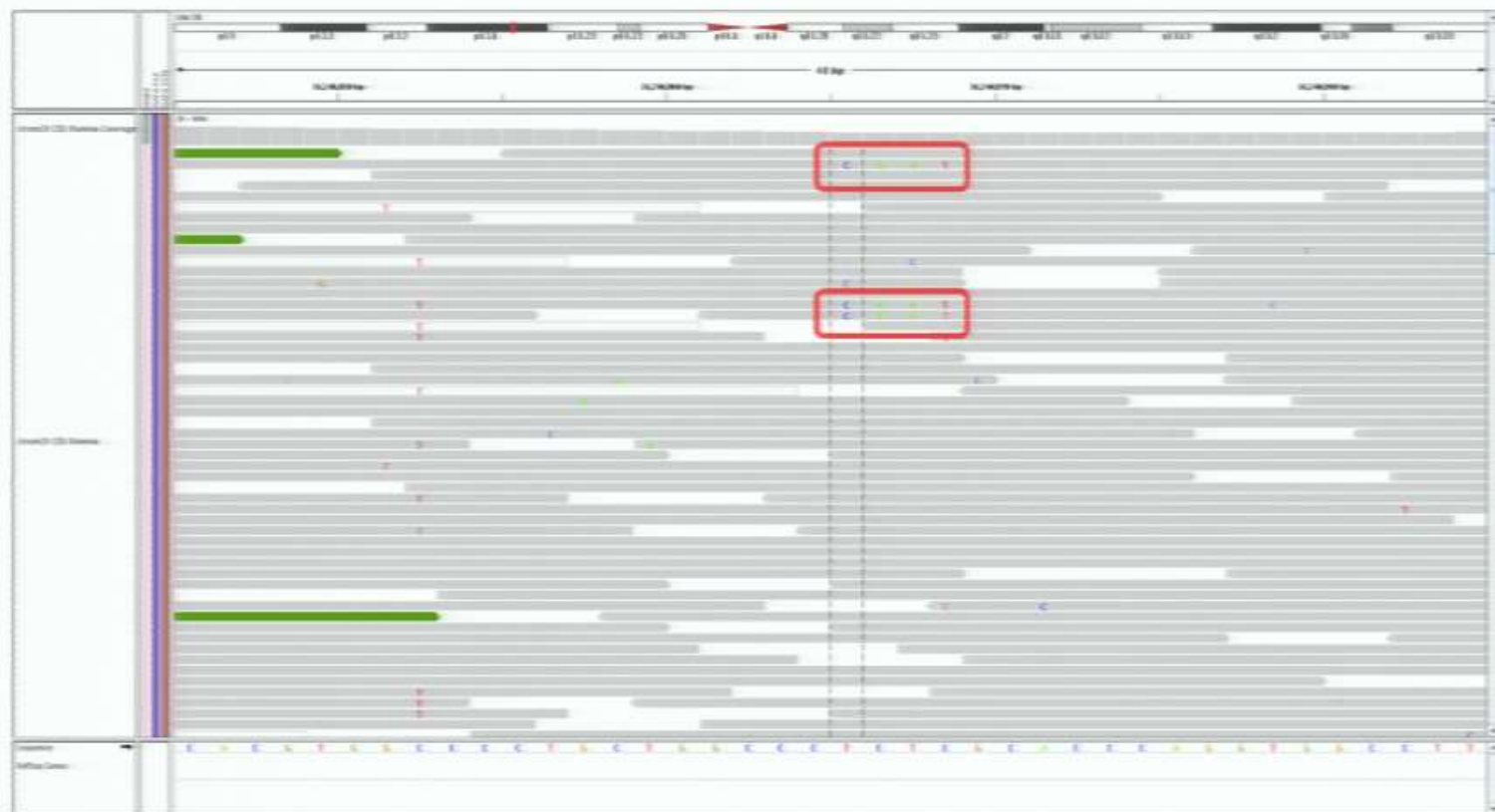
Two yeast strains co-assembled



CCTTGTGTT-----TTAAAT

CCTTGTGTTAGTAATCATATCTTCAGGAAAGCCAAGAACTAATGGCAAA TTAAAT

Example from 1000 Genomes: 4bp MNP



20 16240965 . TCTC CAAT 120 PASS AF=0.31894;NumReads=1107;VarDP=43

NB: only 3 reads in BAM: 43 reads in assembly

Using BWT/FM index for assembly

- The FM-Index of the reads is closely related to an assembly graph
 - It efficiently encodes all 1 base extensions of a k -mer
 - It can compute non-transitive overlap structure of the reads and hence string graph unitigs in $O(N)$ time
 - Also can derive all de Bruijn graphs for k up to read length
- SGA assembler (Jared Simpson)
 - Compression of read BWT vastly reduces memory usage

"Efficient de novo assembly of large genomes using compressed data structures", Simpson and Durbin, 2012

- *Open question: is Myers string graph the same as de Bruijn graph after perfect read threading?*

Final big question

- Can we merge PBWT and BWT assembly?
- Build a model from all existing *primary* data to use efficiently for interpreting new data
 - Use ML techniques, perhaps on probabilistic model capturing genetic structure
- Very large distributed data sets (PB now)
 - The more people, the more information

Is anyone interested to help? Collaborators, postdocs....

- Of course, there are lots of other big data machine learning problems in genomics

Acknowledgements



Andrew Brown, Milan Malinsky,
Stephan Schiffels, Vladimir Shchur,
Vagheesh Narasimhan, Zhihao Ding,
Yasin Memari, *Jared Simpson*,
Heng Li



Thomas Keane, Sendu Bala,
Petr Danecek, Shane McCarthy,
core sequencing teams





BIG & QUIC: Sparse Inverse Covariance Estimation for a Million Variables

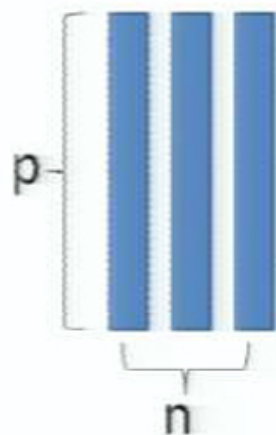
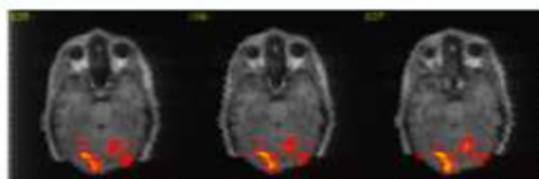
Cho-Jui Hsieh
The University of Texas at Austin

NIPS
Lake Tahoe, Nevada
Dec 8, 2013

Joint work with M. Sustik, I. Dhillon, P. Ravikumar and R. Poldrack

- Goal: Reveal functional connections between regions of the brain.
(Sun et al, 2009; Smith et al, 2011; Varoquaux et al, 2010; Ng et al, 2011)
- $p = 228,483$ voxels.

Input



Output

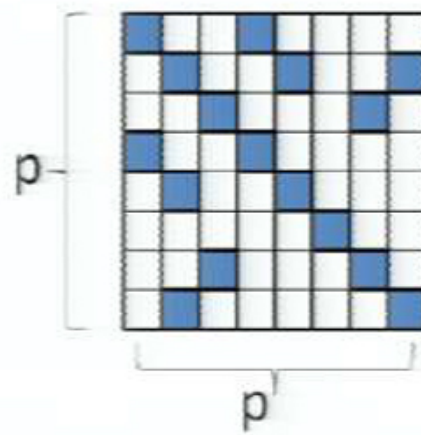
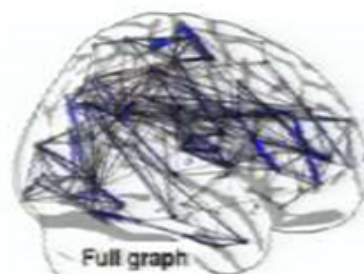


Figure from (Varoquaux et al, 2010)

BIG & QUIC: Sparse Inverse Covariance Estimation for a Million Variables

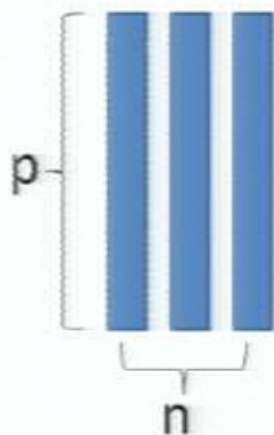
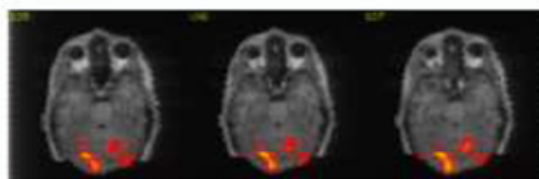
Cho-Jui Hsieh
The University of Texas at Austin

NIPS
Lake Tahoe, Nevada
Dec 8, 2013

Joint work with M. Sustik, I. Dhillon, P. Ravikumar and R. Poldrack

- Goal: Reveal functional connections between regions of the brain.
(Sun et al, 2009; Smith et al, 2011; Varoquaux et al, 2010; Ng et al, 2011)
- $p = 228,483$ voxels.

Input



Output

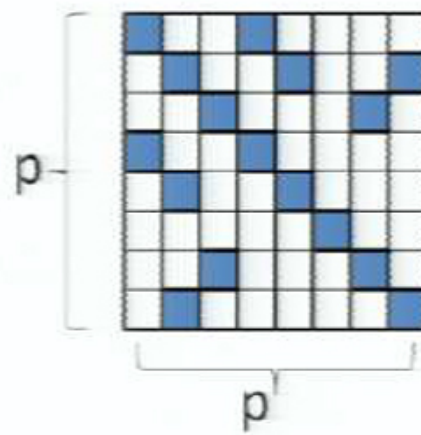
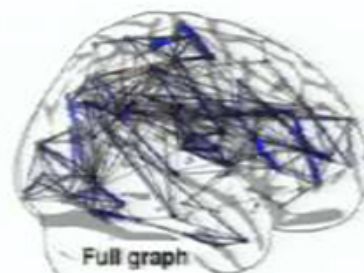


Figure from (Varoquaux et al, 2010)

- Gene regulatory network discovery:
(Schafer & Strimmer 2005; Andrei & Kendzierski 2009; Menendez et al, 2010; Yin and Li, 2011)
- Financial Data Analysis:
 - Model dependencies in multivariate time series (Xuan & Murphy, 2007).
 - Sparse high dimensional models in economics (Fan et al, 2011).
- Social Network Analysis / Web data:
 - Model co-authorship networks (Goldenberg & Moore, 2005).
 - Model item-item similarity for recommender system (Agarwal et al, 2011).
- Climate Data Analysis (Chen et al., 2010).
- Signal Processing (Zhang & Fung, 2013).
- Anomaly Detection (Ide et al, 2009).

Inverse Covariance Estimation

- Given: n i.i.d. samples $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$, $\mathbf{y}_i \in \mathbb{R}^p$, $\mathbf{y}_i \sim \mathcal{N}(\mu, \Sigma)$.
- An example – Chain graph: $y_j = 0.5y_{j-1} + \mathcal{N}(0, 1)$



$$\Sigma = \begin{pmatrix} 1.33 & 0.67 & 0.33 & 0.17 \\ 0.67 & 1.33 & 0.67 & 0.33 \\ 0.33 & 0.67 & 1.33 & 0.67 \\ 0.17 & 0.33 & 0.67 & 1.33 \end{pmatrix}, \quad \Sigma^{-1} = \begin{pmatrix} 1 & -0.5 & 0 & 0 \\ -0.5 & 1.25 & -0.5 & 0 \\ 0 & -0.5 & 1.25 & -0.5 \\ 0 & 0 & -0.5 & 1 \end{pmatrix}$$

- Conditional independence is reflected as zeros in Σ^{-1} :

$\Sigma_{ij}^{-1} = 0 \Leftrightarrow y_i$ and y_j are conditionally independent given other variables.

- Goal: Estimate the **inverse covariance matrix** in the high dimensional setting: $p(\# \text{ variables}) \gg n(\# \text{ samples})$
- Add **ℓ_1 regularization** – a sparse inverse covariance matrix is preferred.
- The ℓ_1 -regularized Maximum Likelihood Estimator:

$$\Sigma^{-1} = \arg \min_{X \succ 0} \left\{ \underbrace{-\log \det X + \text{tr}(SX)}_{\text{negative log likelihood}} + \lambda \|X\|_1 \right\} = \arg \min_{X \succ 0} f(X),$$

where $\|X\|_1 = \sum_{i,j=1}^n |X_{ij}|$.

- Block coordinate ascent (Banerjee et al, 2007), Graphical Lasso (Friedman et al, 2007).
- VSM, PSM, SINCO, IPM, PQN, ALM (2008-2010).
ALM solves $p = 1000$ in 300 secs.
- QUIC: Newton type method (Hsieh et al, 2011)
Solves $p = 1000$ in 10 secs, $p = 10,000$ in half hour.
- All the above methods require $O(p^2)$ memory, cannot solve problems with $p > 30,000$.
- Need for scalability: FMRI dataset has more than 220,000 variables

- Block coordinate ascent (Banerjee et al, 2007), Graphical Lasso (Friedman et al, 2007).
- VSM, PSM, SINCO, IPM, PQN, ALM (2008-2010).
ALM solves $p = 1000$ in 300 secs.
- QUIC: Newton type method (Hsieh et al, 2011)
Solves $p = 1000$ in 10 secs, $p = 10,000$ in half hour.
- All the above methods require $O(p^2)$ memory, cannot solve problems with $p > 30,000$.
- Need for scalability: FMRI dataset has more than 220,000 variables
- BIGQUIC (2013):
 $p = 1,000,000$ (1 trillion parameters) in 22.9 hrs with 32 GBytes memory (using a single machine with 32 cores).

- Main Ingredients:
 - ① Second-order Newton-like method (QUIC)
→ quadratic convergence rate.
 - ② Memory-efficient scheme using block coordinate descent (BigQUIC)
→ scale to one million variables.
 - ③ Approximate Hessian computation (BigQUIC)
→ super-linear convergence rate.

- Split smooth and non-smooth terms: $f(X) = g(X) + h(X)$, where

$$g(X) = -\log \det X + \text{tr}(SX) \text{ and } h(X) = \lambda \|X\|_1.$$

- Form quadratic approximation for $g(X_t + \Delta)$:

$$\begin{aligned} \bar{g}_{X_t}(\Delta) = & \text{tr}((S - W_t)\Delta) + (1/2) \text{vec}(\Delta)^T (W_t \otimes W_t) \text{vec}(\Delta) \\ & - \log \det X_t + \text{tr}(SX_t), \end{aligned}$$

where $W_t = (X_t)^{-1} = \frac{\partial}{\partial X} \log \det(X) |_{X=X_t}$.

- Define the generalized Newton direction:

$$D_t = \arg \min_{\Delta} \bar{g}_{X_t}(\Delta) + \lambda \|X_t + \Delta\|_1.$$

- Solve by coordinate descent (Hsieh et al, 2011) or other methods (Olsen et al, 2012).

- Use coordinate descent to solve:

$$\arg \min_D \{ \bar{g}_X(D) + \lambda \|X + D\|_1 \}.$$

- Closed form solution for each coordinate descent update:

$$D_{ij} \leftarrow -c + S(c - b/a, \lambda/a),$$

where $S(z, r) = \text{sign}(z) \max\{|z| - r, 0\}$ is the soft-thresholding function, $a = W_{ij}^2 + W_{ii}W_{jj}$, $b = S_{ij} - W_{ij} + \mathbf{w}_i^T D \mathbf{w}_j$, and $c = X_{ij} + D_{ij}$.

- The main cost is in computing $\mathbf{w}_i^T D \mathbf{w}_j$,

where $\mathbf{w}_i, \mathbf{w}_j$ are i -th and j -th columns of $W = X^{-1}$.

QUIC: QUadratic approximation for sparse Inverse Covariance estimation

Input: Empirical covariance matrix S , scalar λ , initial X_0 .

For $t = 0, 1, \dots$

- 1 Variable selection: select a *free* set of $m \ll p^2$ variables.
- 2 Use coordinate descent to find descent direction:
 $D_t = \arg \min_{\Delta} \bar{f}_{X_t}(X_t + \Delta)$ over set of free variables, (A *Lasso* problem.)
- 3 Line Search: use an *Armijo*-rule based step-size selection to get α s.t.
 $X_{t+1} = X_t + \alpha D_t$ is
 - positive definite,
 - satisfies a sufficient decrease condition $f(X_t + \alpha D_t) \leq f(X_t) + \alpha \sigma \Delta_t$.
 (Cholesky factorization of $X_t + \alpha D_t$)

Consider the case that $p \approx 1\text{million}$, $m = \|X_t\|_0 \approx 50\text{million}$.

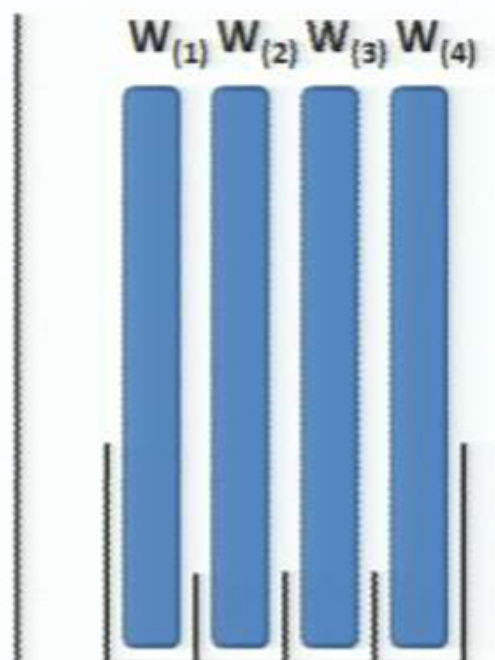
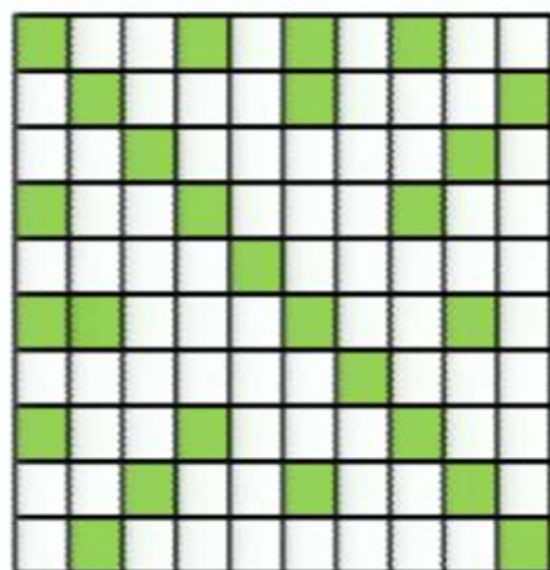
- Coordinate descent requires X_t and $W = X_t^{-1}$,
 - needs $O(p^2)$ storage
 - needs $O(mp)$ computation per sweep, where $m = \|X_t\|_0$
- Line search (compute determinant using Cholesky factorization).
 - needs $O(p^2)$ storage
 - needs $O(p^3)$ computation

Coordinate Updates with Memory Cache

- Assume we can store M columns of W in memory.
- Coordinate descent update (i, j) : compute $\mathbf{w}_i^T D \mathbf{w}_j$.
- If $\mathbf{w}_i, \mathbf{w}_j$ are not in memory: recompute by CG:
 $X \mathbf{w}_i = \mathbf{e}_i$: $O(T_{CG})$ time.

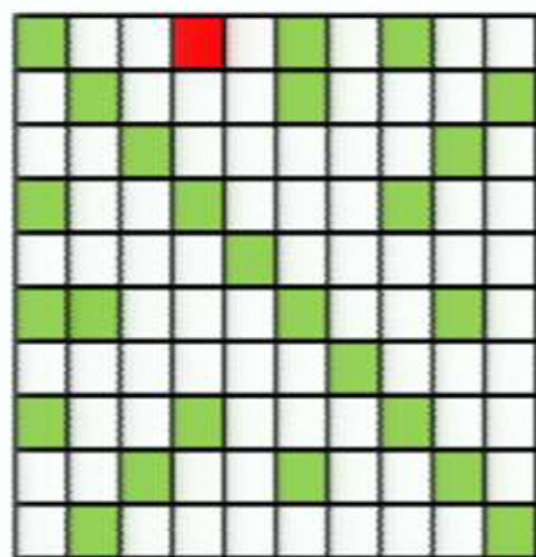
Coordinate Updates with Memory Cache

w_1, w_2, w_3, w_4 stored in memory.

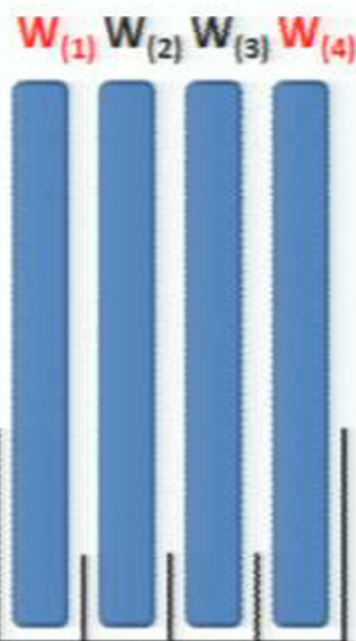


Memory Cache
($M=4$)

Cache hit, do not need to recompute w_i, w_j .

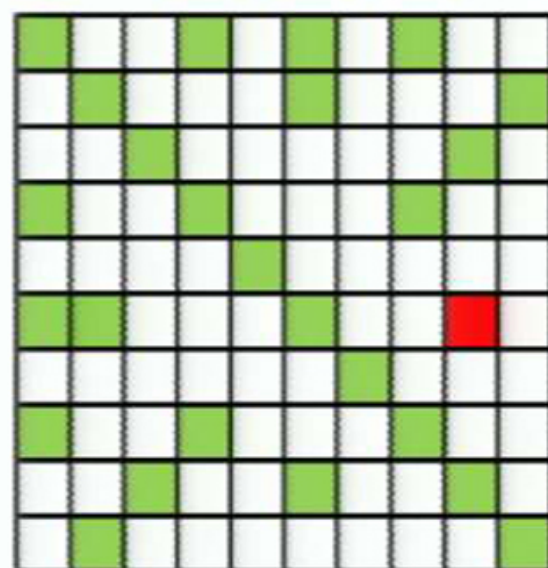


Update (1,4)
Need $w_{(1)}, w_{(4)}$

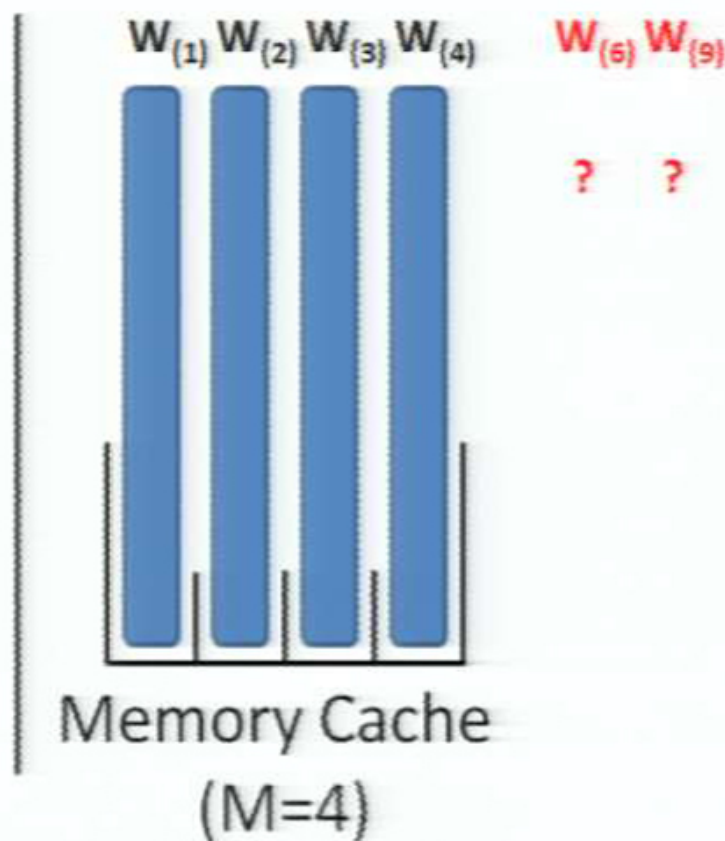


Memory Cache
(M=4)

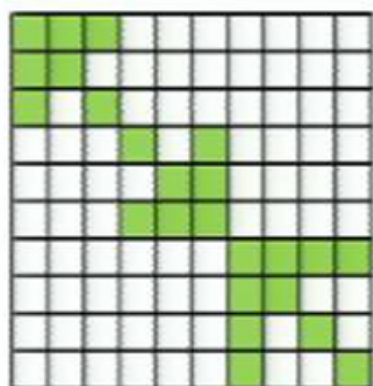
Cache miss, recompute w_i, w_j .



Update (6,9)
Need $w_{(6)}, w_{(9)}$



- Want to find update sequence that **minimizes number of cache misses**: probably NP Hard.
- Our strategy: update variables **block by block**.
- The ideal case: there exists a partition $\{S_1, \dots, S_k\}$ such that all free sets are in diagonal blocks:



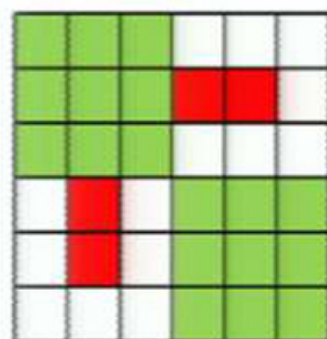
Free Set

- Only requires p column evaluations.

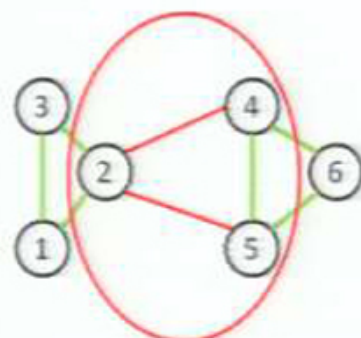
- If the block partition is not perfect:
extra column computations can be characterized by **boundary nodes**.
- Given a partition $\{S_1, \dots, S_k\}$, we define boundary nodes as

$$B(S_q) \equiv \{j \mid j \in S_q \text{ and } \exists i \in S_z, z \neq q \text{ s.t. } F_{ij} = 1\},$$

where F is adjacency matrix of the free set.



Free Set



boundary
nodes

- The number of columns to be computed in one sweep is

$$p + \sum_q |B(S_q)|.$$

- Can be upper bounded by

$$p + \sum_q |B(S_q)| \leq p + \sum_{z \neq q} \sum_{i \in S_z, j \in S_q} F_{ij}.$$

- Use Graph Clustering (METIS or Graclus) to find the partition.
- Example: on fMRI dataset ($p = 0.228$ million) with 20 blocks,
random partition: need 1.6 million column computations.
graph clustering: need 0.237 million column computations.

- Block co-ordinate descent with clustering,
 - needs $O(p^2) \rightarrow O(m + p^2/k)$ storage
 - needs $O(mp) \rightarrow O(mp)$ computation per sweep, where $m = \|X_t\|_0$
- Line search (compute determinant of a big sparse matrix).
 - needs $O(p^2)$ storage
 - needs $O(p^3)$ computation

- Given sparse matrix $A = X_t + \alpha D$, we need to
 - 1 Check its positive definiteness.
 - 2 Compute $\log \det(A)$.
- Our approach computes $\log \det(A)$ in $O(mp)$ time.
- Cholesky factorization in QUIC requires $O(p^3)$ computation.
- If $A = \begin{pmatrix} a & \mathbf{b}^T \\ \mathbf{b} & C \end{pmatrix}$,
 - $\det(A) = \det(C)(a - \mathbf{b}^T C^{-1} \mathbf{b})$
 - A is positive definite iff C is positive definite and $(a - \mathbf{b}^T C^{-1} \mathbf{b}) > 0$.
- C is sparse, so can compute $C^{-1} \mathbf{b}$ using Conjugate Gradient (CG).
- Time complexity: $T_{CG} = O(mT)$, where T is number of CG iterations.

- Block co-ordinate descent with clustering,
 - needs $O(p^2) \rightarrow O(m + p^2/k)$ storage
 - needs $O(mp) \rightarrow O(mp)$ computation per sweep, where $m = \|X_t\|_0$
- Line search (compute determinant of a big sparse matrix).
 - needs $O(p^2) \rightarrow O(p)$ storage
 - needs $O(p^3) \rightarrow O(mp)$ computation

BIGQUIC

Input: Samples Y , scalar λ , initial X_0 .

For $t = 0, 1, \dots$

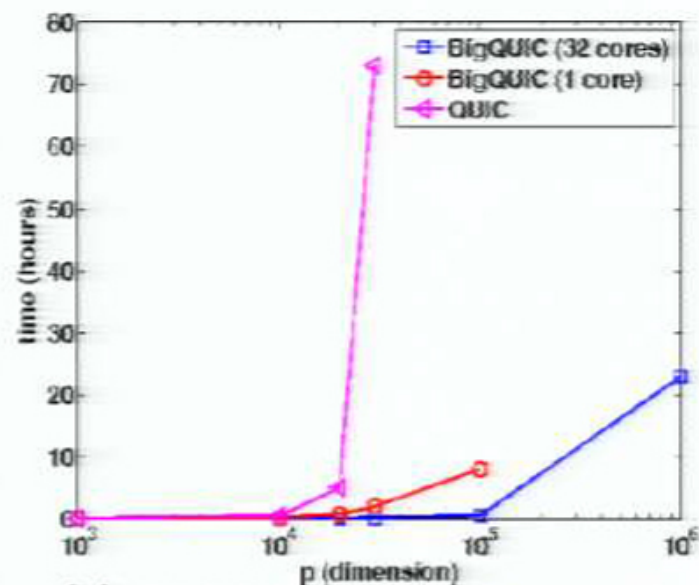
- ① Variable selection: select a *free* set of $m \ll p^2$ variables.
- ② **Construct a partition by clustering.**
- ③ Run **block coordinate descent** to find descent direction:
 $D_t = \arg \min_{\Delta} \bar{f}_{X_t}(X_t + \Delta)$ over set of free variables.
- ④ Line Search: use an *Armijo*-rule based step-size selection to get α s.t.
 $X_{t+1} = X_t + \alpha D_t$ is
 - positive definite,
 - satisfies a sufficient decrease condition $f(X_t + \alpha D_t) \leq f(X_t) + \alpha \sigma \Delta_t$.

(Schur complement with conjugate gradient method.)

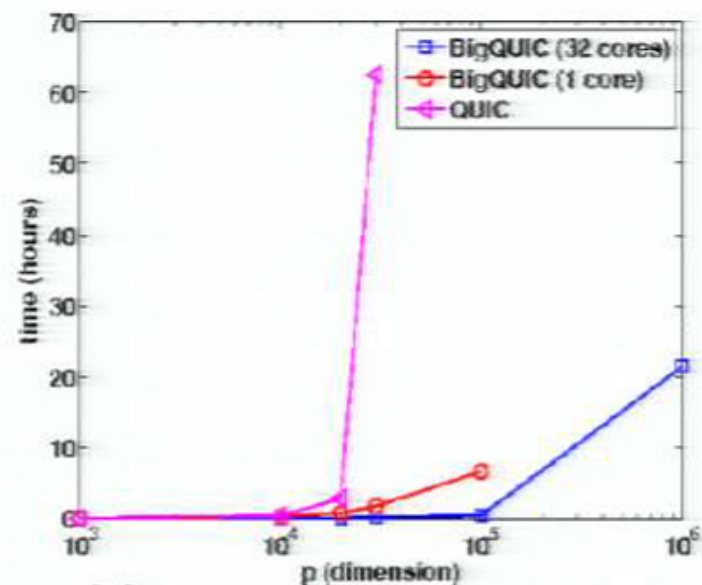
- Recall $W = X^{-1}$.
- When each \mathbf{w}_i is computed by CG ($X\mathbf{w}_i = \mathbf{e}_i$):
 - The gradient $\nabla_{ij}g(X) = S_{ij} - W_{ij}$ on free set can be computed **once** and stored in memory.
 - Hessian ($\mathbf{w}_i^T D \mathbf{w}_j$ in coordinate updates) needs to be **repeatedly computed**.
- To reduce the time overhead, Hessian should be computed **approximately**.
- Theorem:** the convergence rate is quadratic if $\|X\hat{\mathbf{w}}_i - \mathbf{e}_i\| = O(\|\nabla^S f(X_t)\|)$, where

$$\nabla_{ij}^S f(X) = \begin{cases} \nabla_{ij}g(X) + \text{sign}(X_{ij})\lambda & \text{if } X_{ij} \neq 0, \\ \text{sign}(\nabla_{ij}g(X)) \max(|\nabla_{ij}g(X)| - \lambda, 0) & \text{if } X_{ij} = 0. \end{cases}$$

Experimental results (scalability)



(a) Scalability on random graph



(b) Scalability on chain graph

Figure: BIGQUIC can solve one million dimensional problems.

BIGQUIC is faster even for medium size problems.

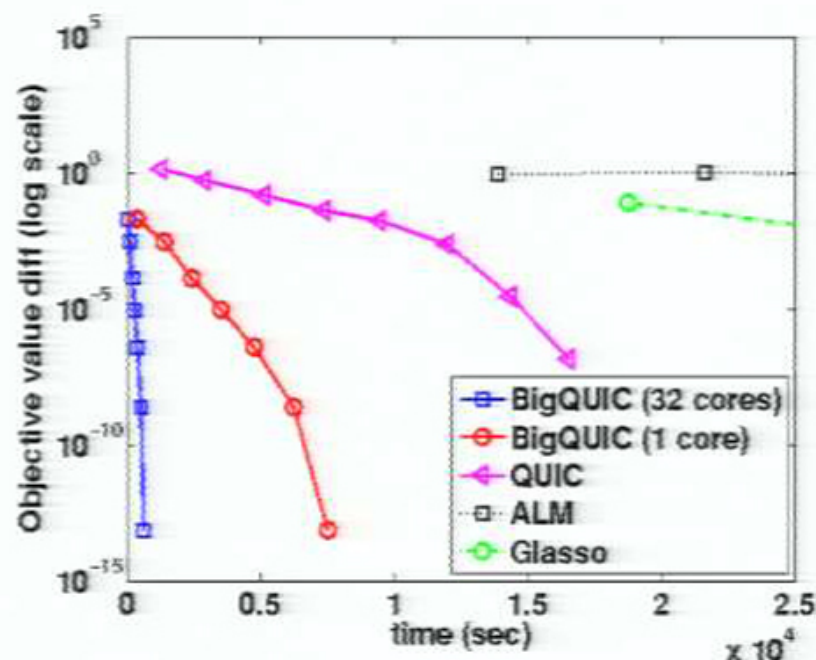
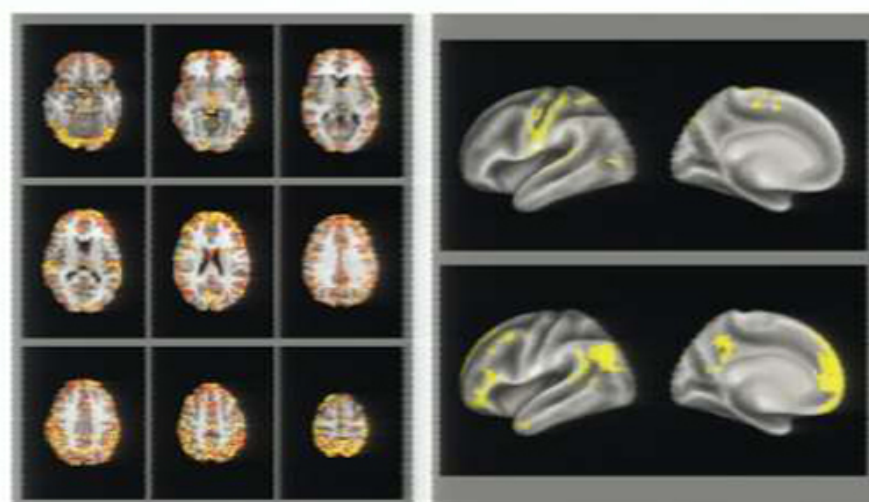


Figure: Comparison on FMRI data with a $p = 20000$ subset (maximum dimension that previous methods can handle).

- 228,483 voxels, 518 time points.
- $\lambda = 0.6 \implies$ average degree 8, BIGQUIC took 5 hours.
 $\lambda = 0.5 \implies$ average degree 38, BIGQUIC took 21 hours.
- Findings:
 - Voxels with large degree were generally found in the gray matter.
 - Can detect meaningful brain modules by modularity clustering.



- BIGQUIC: Memory efficient quadratic approximation method for sparse inverse covariance estimation.
- Our contributions:
 - Computing Newton direction:
 - Coordinate descent \rightarrow **block coordinate descent with clustering.**
 - Memory complexity: $O(p^2) \rightarrow O(m + p^2/k)$.
 - Time complexity: $O(mp) \rightarrow O(mp)$.
 - Line search (computing determinant of a big sparse matrix)
 - Cholesky factorization \rightarrow **Schur complement with conjugate gradient method.**
 - Memory complexity: $O(p^2) \rightarrow O(p)$.
 - Time complexity: $O(p^3) \rightarrow O(mp)$.
 - **Inexact Hessian computation with super-linear convergence.**

- [1] C. J. Hsieh, M. Sustik, I. S. Dhillon, P. Ravikumar, and R. Poldrack. *BIG & QUIC: Sparse inverse covariance estimation for a million variables*. NIPS (oral presentation), 2013.
- [2] C. J. Hsieh, M. Sustik, I. S. Dhillon, and P. Ravikumar. *Sparse Inverse Covariance Matrix Estimation using Quadratic Approximation*. NIPS, 2011.
- [3] C. J. Hsieh, I. S. Dhillon, P. Ravikumar, A. Banerjee. *A Divide-and-Conquer Procedure for Sparse Inverse Covariance Estimation*. NIPS, 2012.
- [4] P. A. Olsen, F. Oztoprak, J. Nocedal, and S. J. Rennie. *Newton-Like Methods for Sparse Inverse Covariance Estimation*. Optimization Online, 2012.
- [5] O. Banerjee, L. El Ghaoui, and A. d'Aspremont. *Model Selection Through Sparse Maximum Likelihood Estimation for Multivariate Gaussian or Binary Data*. JMLR, 2008.
- [6] J. Friedman, T. Hastie, and R. Tibshirani. *Sparse inverse covariance estimation with the graphical lasso*. Biostatistics, 2008.
- [7] L. Li and K.-C. Toh. *An inexact interior point method for l_1 -regularized sparse covariance selection*. Mathematical Programming Computation, 2010.
- [8] K. Scheinberg, S. Ma, and D. Glodfarb. *Sparse inverse covariance selection via alternating linearization methods*. NIPS, 2010.
- [9] K. Scheinberg and I. Rish. *Learning sparse Gaussian Markov networks using a greedy coordinate ascent approach*. Machine Learning and Knowledge Discovery in Databases, 2010.