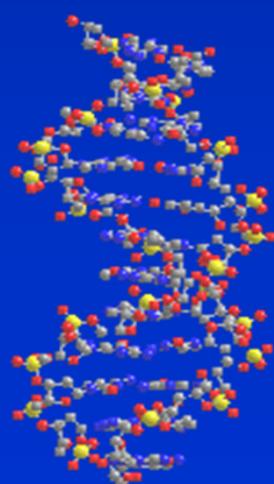


Genome wide association and Genomic Selection in the era of Genome sequencing



Course overview

- Day 1
 - Quantitative traits
 - Linkage disequilibrium
 - Genome wide association studies
- Day 2 and 3
 - Genomic prediction - BLUP and GBLUP
 - Genomic prediction – Bayesian methods
- Day 4
 - Validation of genomic predictions
 - Optimal breeding program design with genomic selection
- Day 5
 - Imputation and whole genome sequencing for genomic selection

Course overview

- Day 1
 - Quantitative traits
 - Genome wide association studies
- Day 2 and 3
 - Genomic prediction - BLUP
 - Genomic prediction – Bayesian methods
- Day 4
 - Imputation and whole genome sequencing for genomic selection

Imputation

- Why impute?
- Approaches for imputation
- Factors affecting accuracy of imputation
- How can imputation give you more power?

Why impute?

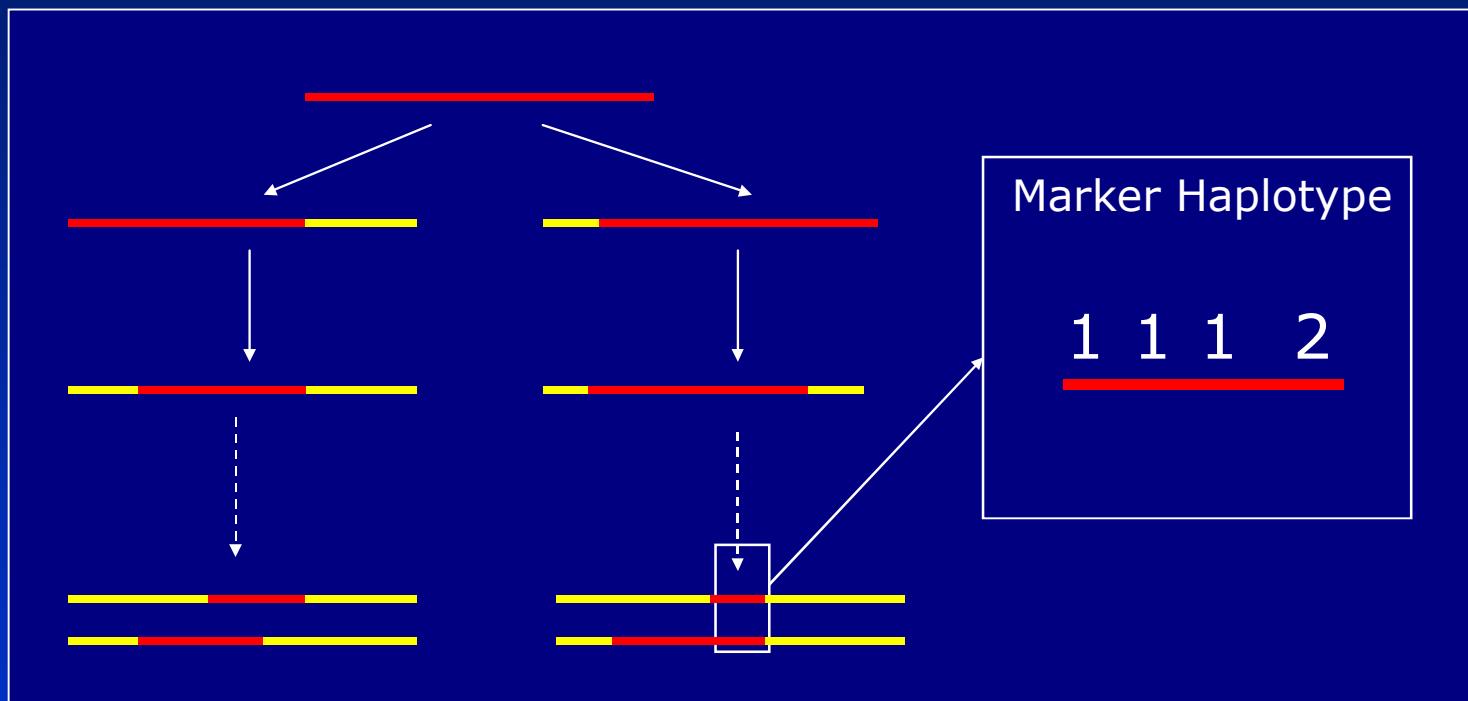
- Fill in missing genotypes from the lab
- Merge data sets with genotypes on different arrays
 - Eg. Affy and Illumina data
- Impute from low density to high density
 - 7K-> 50K (save \$\$\$)
 - 50K->800K
 - capture power of higher density?
 - Better persistence of accuracy
- Sequence expensive, can we impute to full sequence data?

Core concept

- Identity by state (IBS)
 - A pair of individuals have the same allele at a locus
- Identity by descent (IBD)
 - A pair of individuals have the same alleles at a locus and it traces to a common ancestor
- Imputation methods determine whether a chromosome segment is IBD

Causes of LD

- A chunk of ancestral chromosome is conserved in the current population



Core concept 2

- Any individuals in a population may share a proportion of their genome identical by descent (IBD)
 - IBD segments are the same and have originated in a common ancestor
- The closer the relationship the longer the IBD segments
 - Pedigree relationships

Several methods for imputation

- Two main categories:
 - Family based
 - Population based
 - Or combination of the two
 - Some of the most effective are Beagle (Browning and Browning, 2009), MACH (Li et al., 2010), Impute2 (Howie et al., 2009), AlphaPhase (Hickey et al 2011)

Several methods for imputation

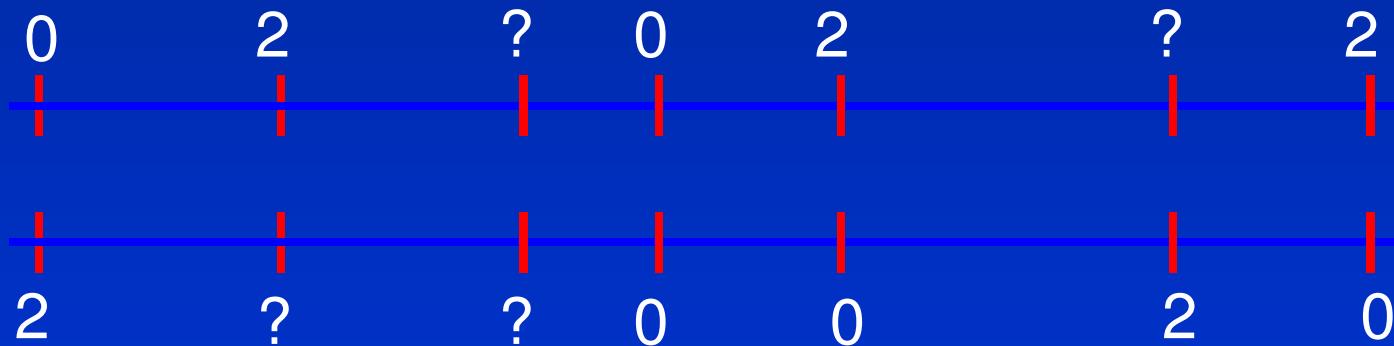
- Two main categories:
 - Family based
 - Population based
 - Or combination of the two
 - Some of the most effective are Beagle (Browning and Browning, 2009), MACH (Li et al., 2010), Impute2 (Howie et al., 2009), AlphaPhase (Hickey et al 2011)

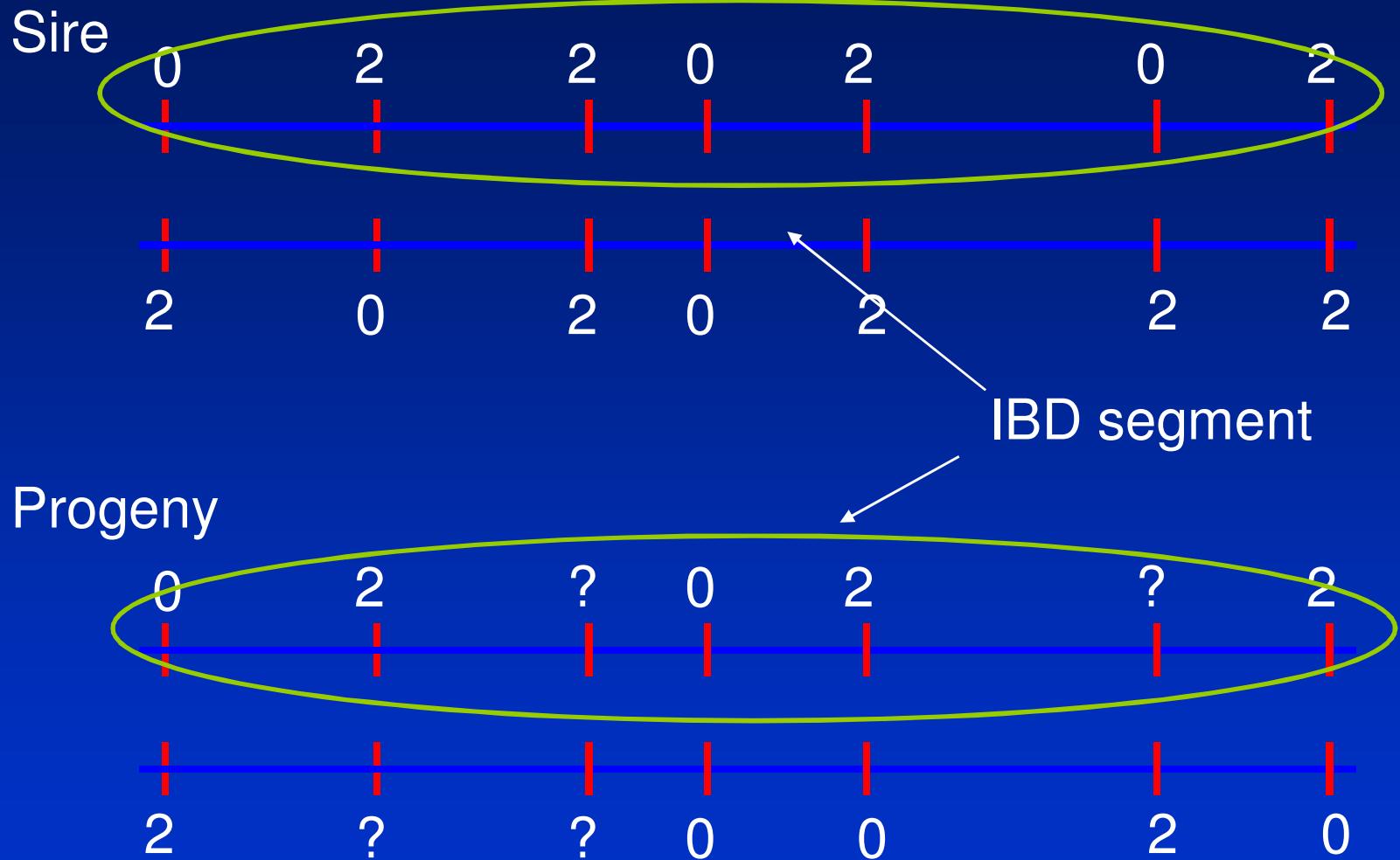
Finding an IBD segment

Sire

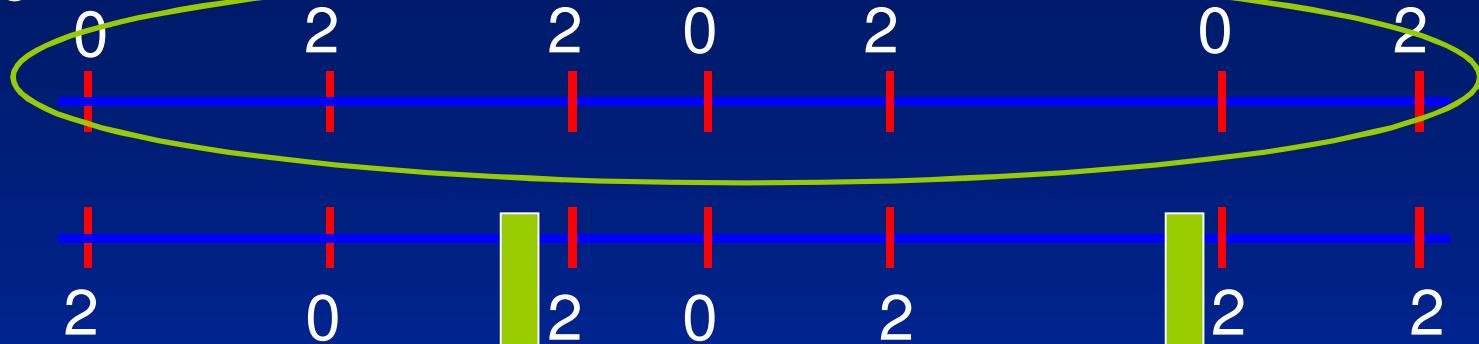


Progeny

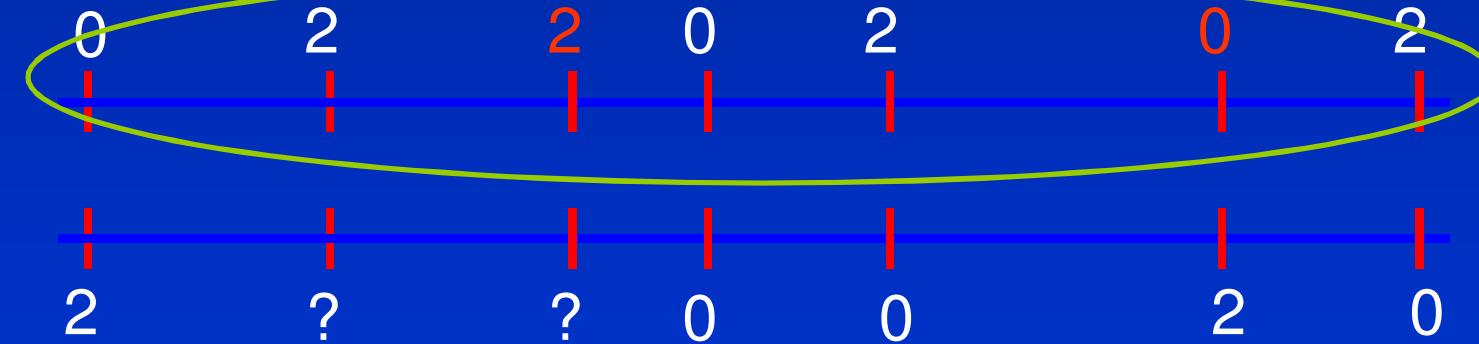




Sire



Progeny



Several methods for imputation

- Two main categories:
 - Family based
 - Population based (exploits LD)
 - Or combination of the two
 - Some of the most effective are Beagle (Browning and Browning, 2009), MACH (Li et al., 2010), Impute2 (Howie et al., 2009), AlphaPhase (Hickey et al 2011)

Population based imputation

- Hidden Markov Models
 - Has “hidden states”
 - For target individuals these are “map” of reference haplotypes that have been inherited
 - Imputation problem is to derive genotype probabilities given hidden states, sparse genotypes, recombination rates, other population parameters

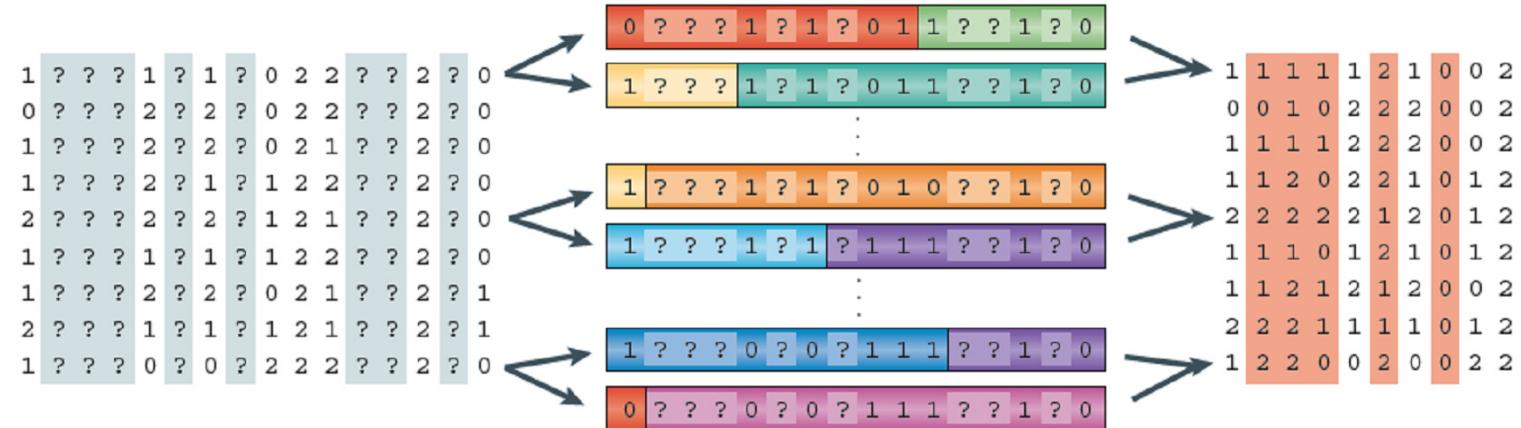
$$P(G_i|H, \theta, \rho) = \sum_s P(G_i|s, \theta)P(s|H, \rho)$$

Population based imputation

Reference population

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

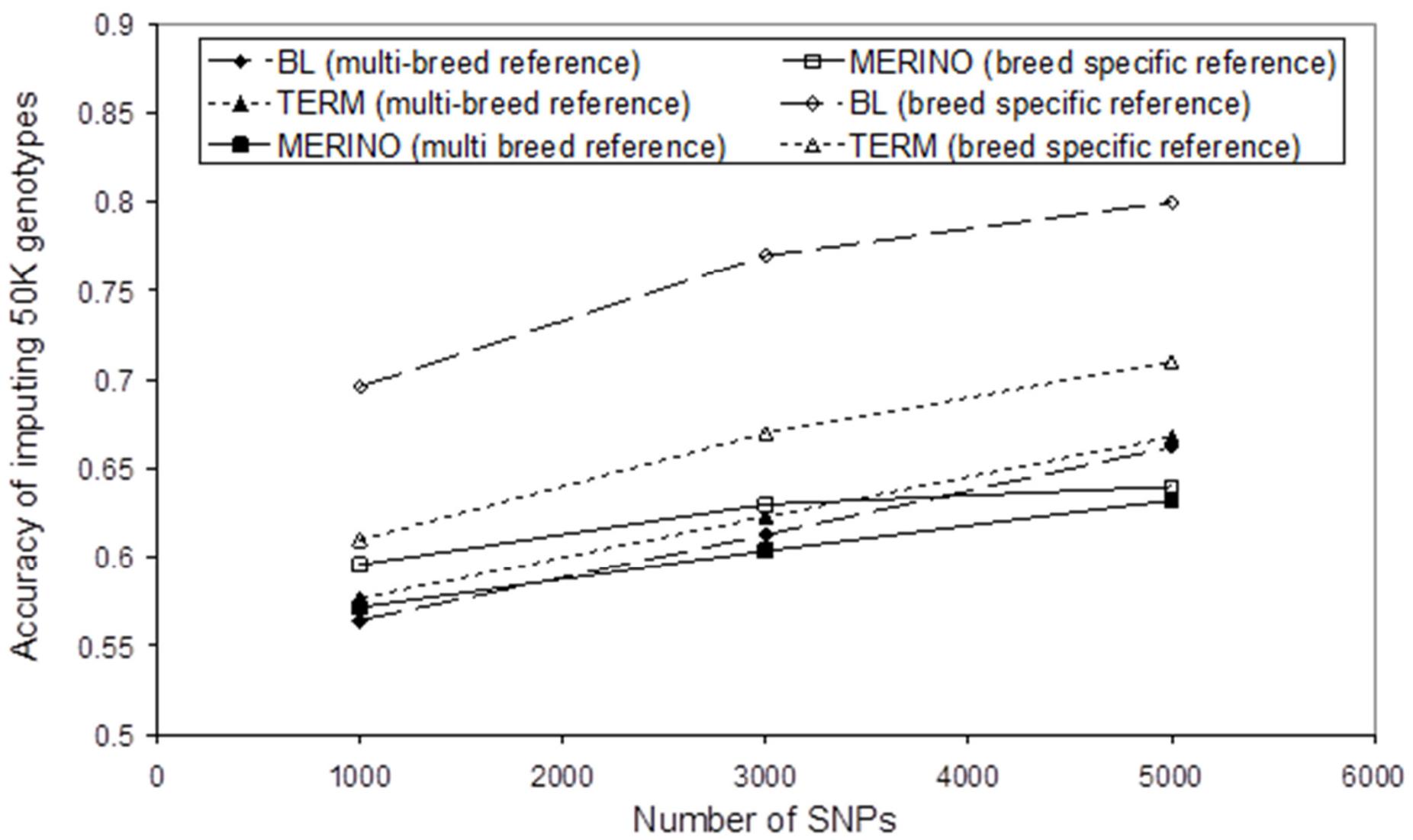
Target population



Imputation accuracy

- Depends on
 - Size of reference set
 - bigger the better!
 - Density of markers
 - extent of LD, effective population size
 - Frequency of SNP alleles
 - Genetic relationship to reference

Imputation accuracy sheep



Imputation accuracy

- Density of markers (extent of LD)
 - In Holstein Dairy cattle
 - 3K -> 50K accuracy 0.93
 - 7K -> 50K accuracy 0.98

Illumina Bovine HD array

- We genotyped
 - 898 Holstein heifers
 - 47 Holstein Key ancestor bulls
- After (stringent) QC **634,307** SNPs

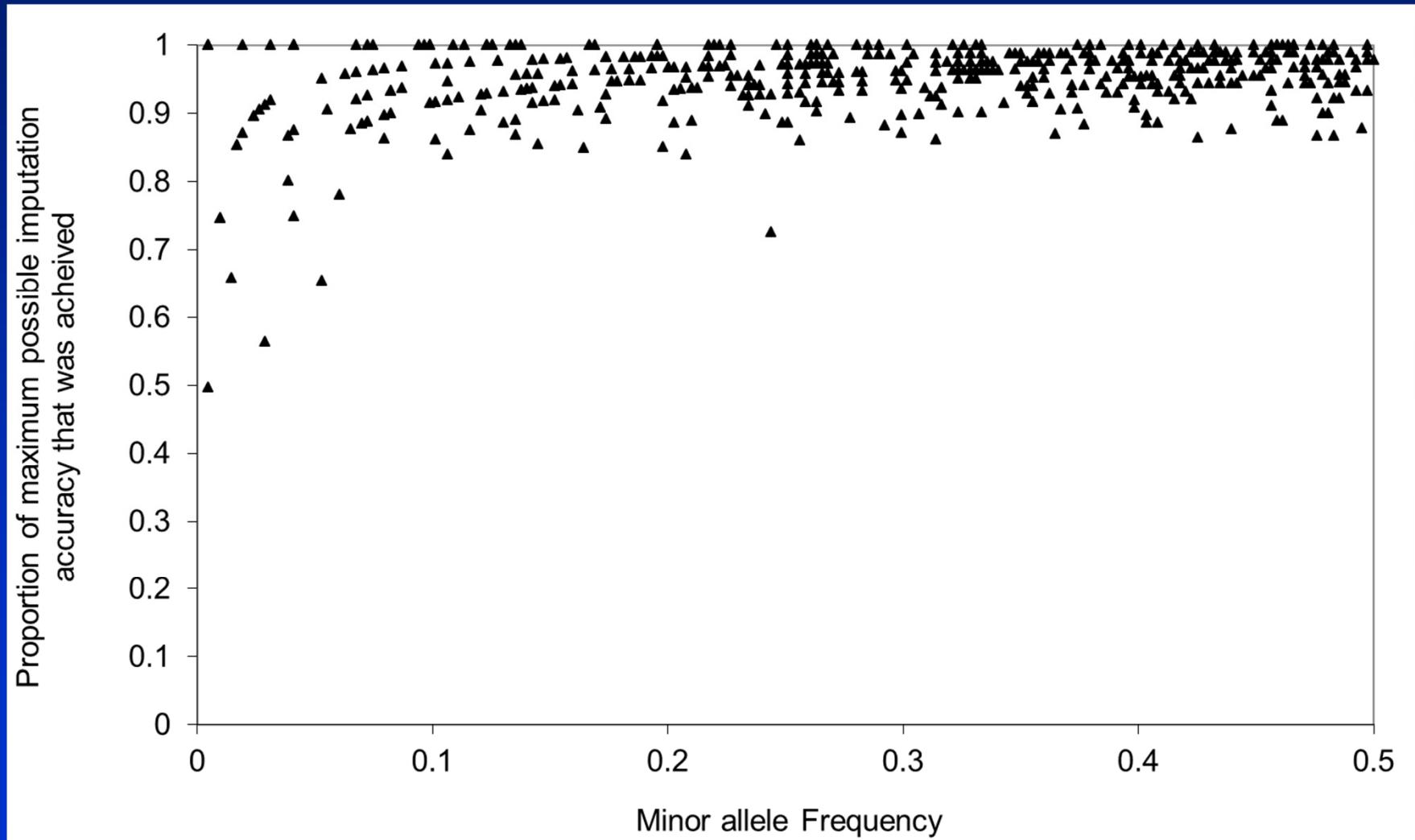
Imputation 50K -> 800K

- Holsteins

	Cross validation	% Correct
Heifers only	1	96.7%
	2	96.7%
	Average	96.7%
Heifers using key ancestors	1	97.8%
	2	97.7%
	Average	97.7%

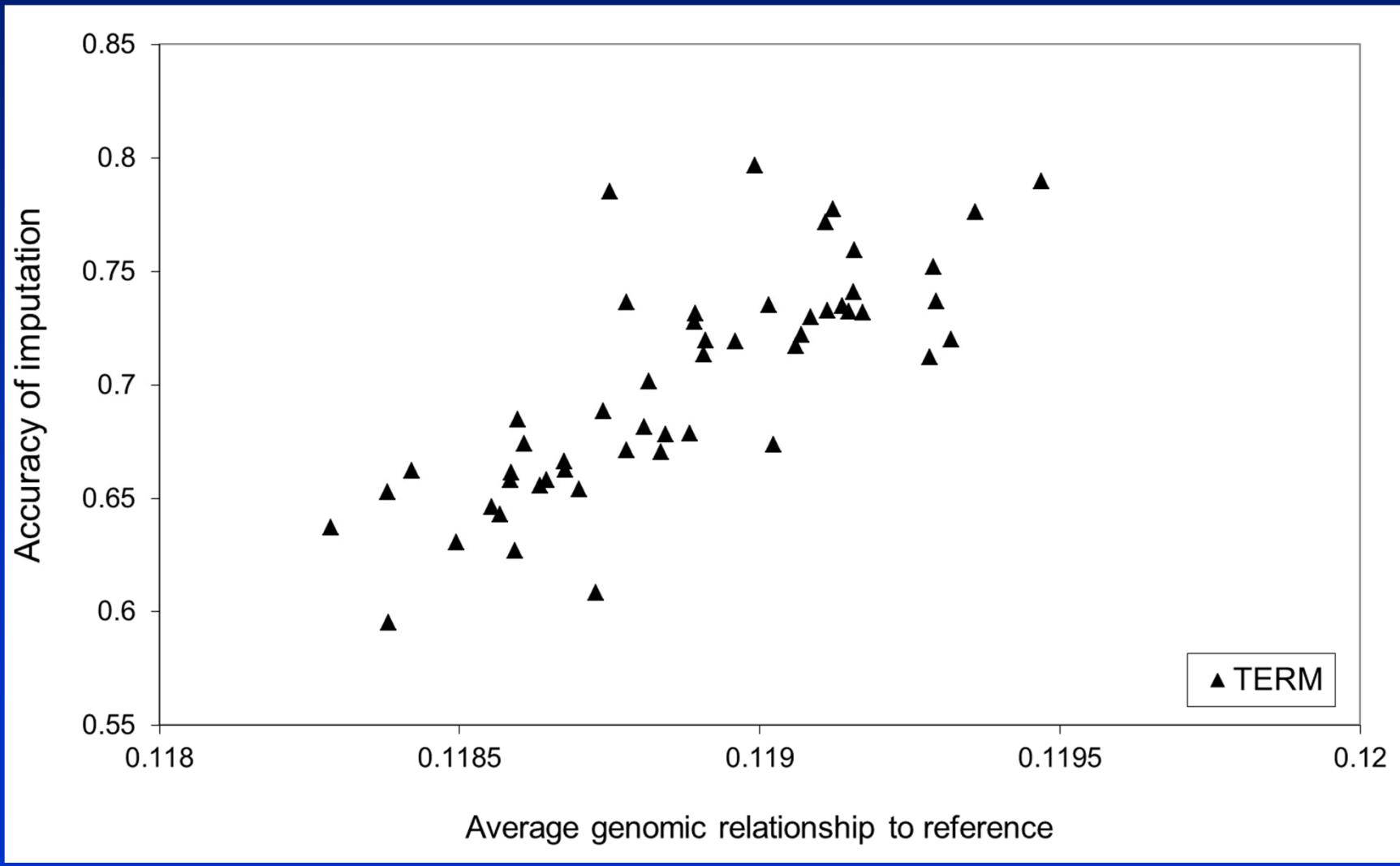
Imputation accuracy

- Rare alleles?



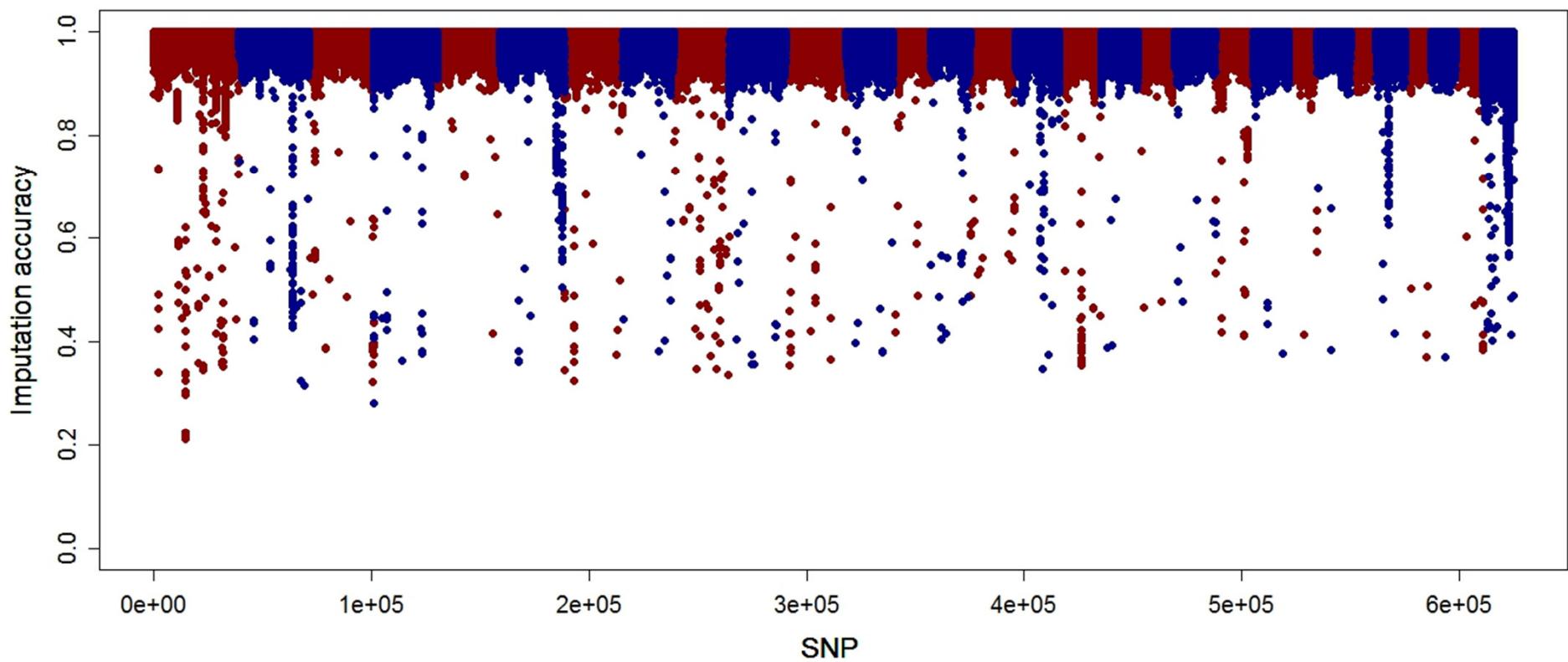
Imputation accuracy

- Relationship to reference?



Imputation of full sequence data

- Effect of map errors?



Why more power with imputation

- High accuracies of imputation demonstrate that we can infer haplotypes of animal genotyped with e.g. 3K accurately
- But potentially large number of haplotypes
- With imputed data can test single snp, only use 1 degree of freedom, rather than number of haplotypes

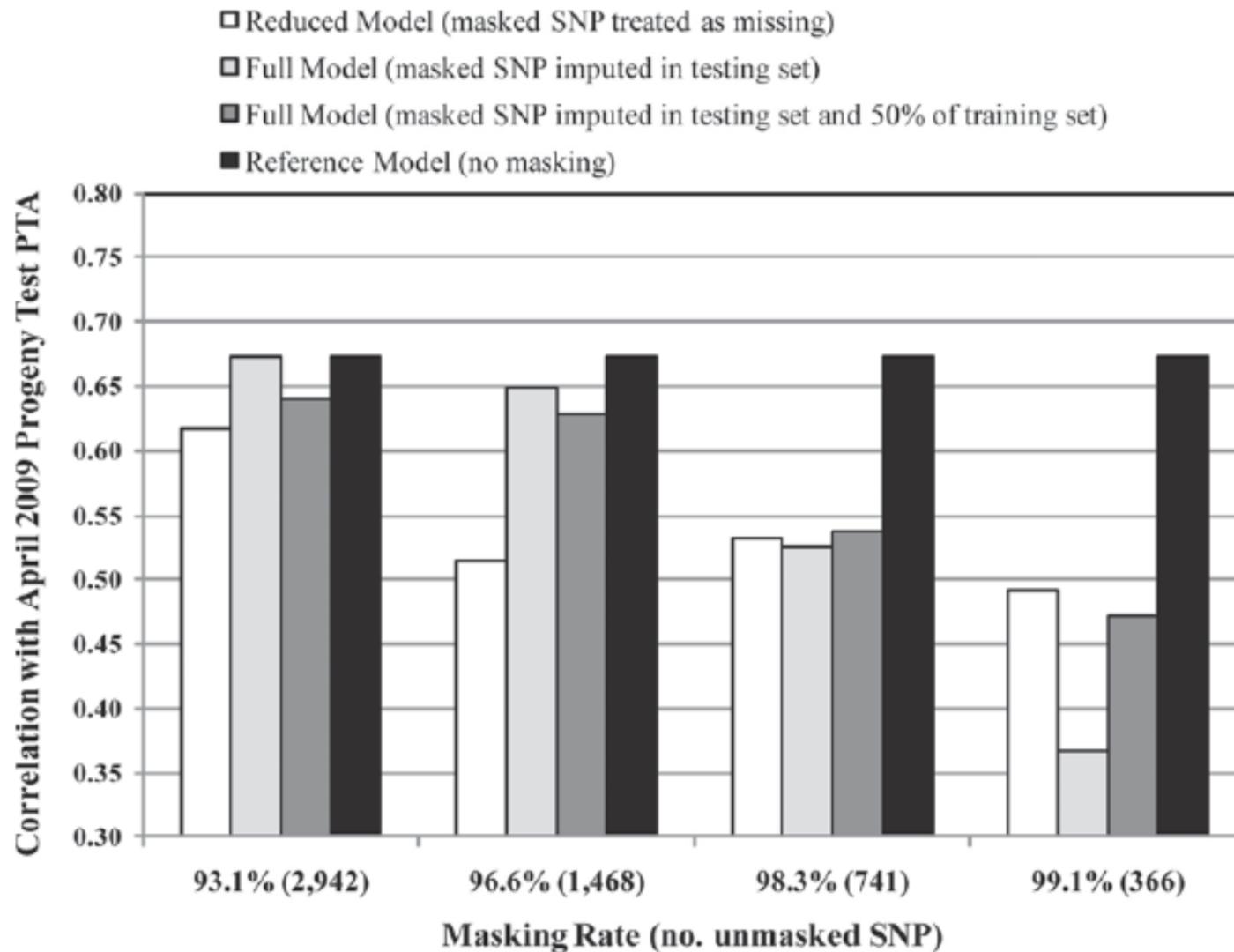


Figure 2. Correlations between predicted direct genomic values for milk yield and corresponding April 2009 progeny-test PTA using full or reduced models with 42,552 or 366, 741, 1,468, or 2,942 single SNP covariates, respectively, with or without imputation of masked genotypes for bulls in the testing set or bulls in the testing set and a randomly chosen 50% of bulls in the training set. The bars denoted as “reference” correspond to correlations from a full model in which all 42,552 SNP genotypes were left as unmasked in both the training and testing sets.

Using sequence data in genomic selection and GWAS

- Motivation
- Characteristics of sequence data
- Which individuals to sequence?
- Imputation of full sequence data
- Methods for genomic prediction with full sequence data
- Examples
 - GWAS in Rice, Cattle

Using sequence data in genomic selection and GWAS

- Motivation
 - Genome wide association study
 - Straight to causative mutation
 - Genomic selection (all hypotheses!)
 - No longer have to rely on LD, causative mutation actually in data set
 - Higher accuracy of prediction?
 - Better prediction across breeds?
 - Assumes same QTL segregating in both breeds
 - No longer have to rely on SNP-QTL associations holding across breeds
 - Better persistence of accuracy across generations

Using sequence data in genomic selection and GWAS

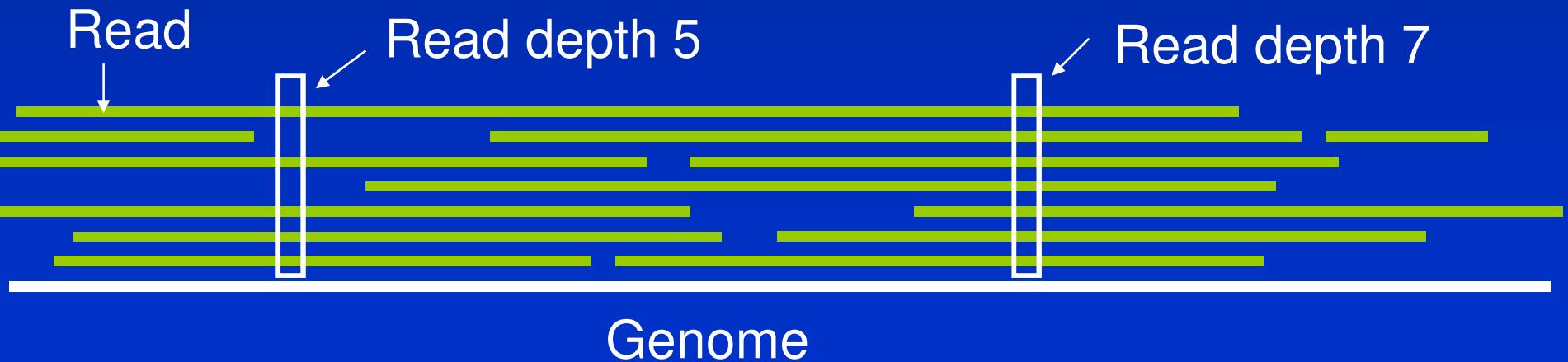
- Motivation
- Characteristics of sequence data
- Which individuals to sequence?
- Imputation of full sequence
- Methods for genomic prediction with full sequence data
- Examples
 - GWAS in Rice, Cattle

Sequence data

- Generates reads of DNA approx. 100 base pair (bp) length
- Reads are aligned to a reference genome
 - Or they could be assembled *de novo*
 - Assigns each read a location on genome
- Reads have an error rate!
 - One error per read
- Information is base pair (ACTG) + Quality score for each base
 - PHRED score = $-10 \cdot \log_{10}(\text{error rate})$
 - 0.01 error rate = Q20
 - 0.001 error rate = Q30
 - 0.0001 error rate = Q40

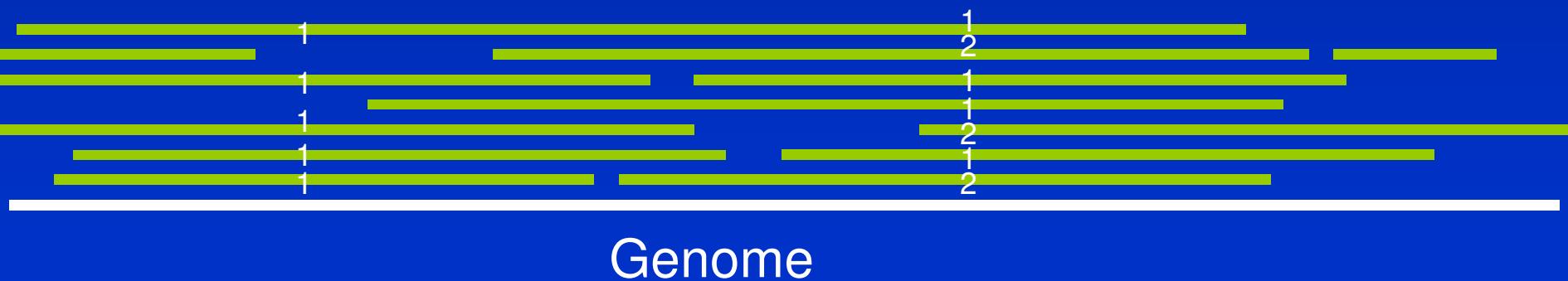
Read depth

- Each sequenced animal is aligned separately to reference
 - .bam files are created
- Read depth or fold coverage



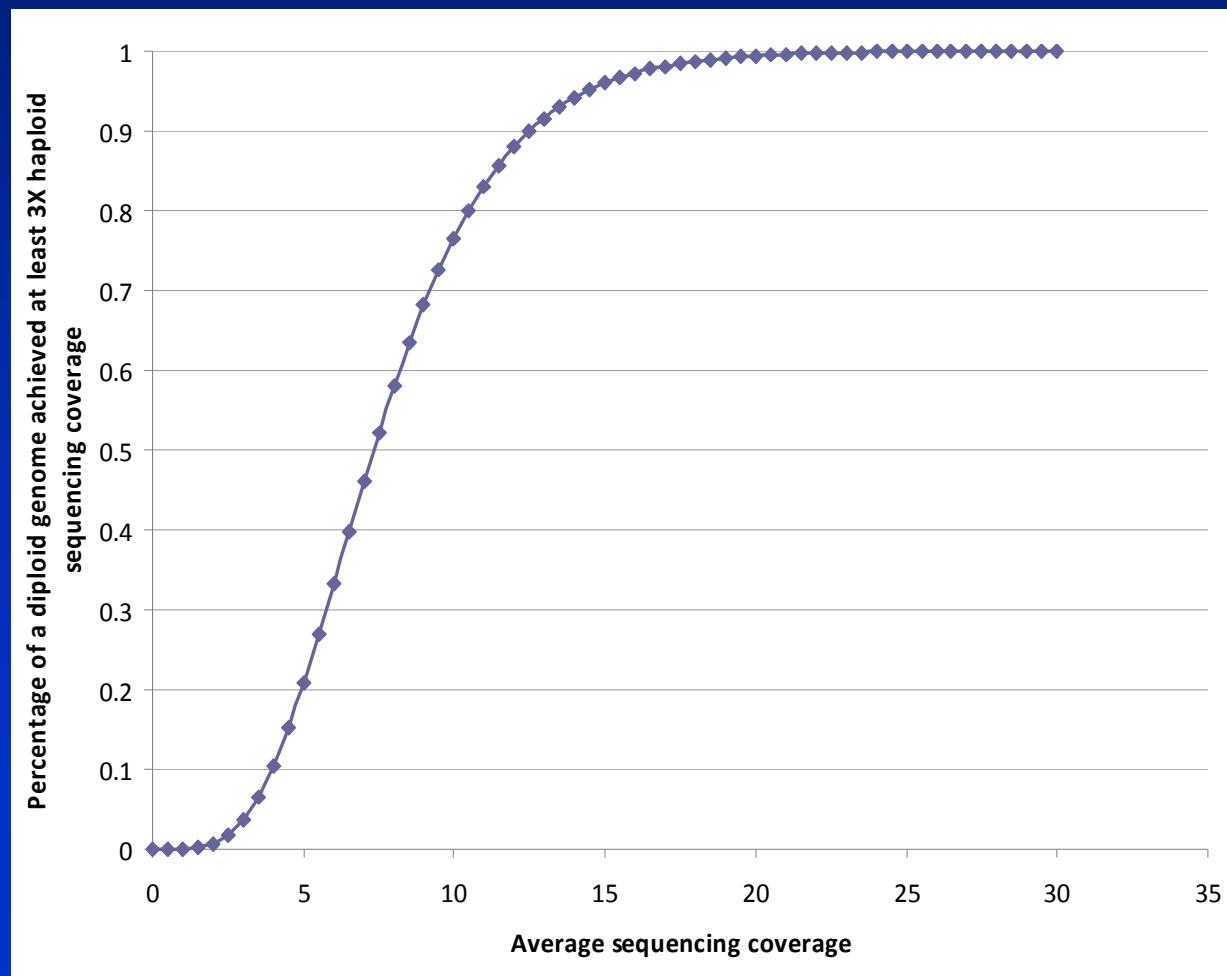
Importance of read depth

- Consider a heterozygous locus (animal carries 2 different alleles)
 - 50/50 chance of observing each allele in every read
- If read depth is low, it is possible to not observe an allele and therefore call a heterozygous locus homozygous
 - Read depth 5 → $0.5^5 = 0.03125$



What read depth is sufficient?

- Proportion of genome achieving at least 6x diploid coverage
- 12.5x achieves 90% in simulation below (Shen et al. 2010, Suppl. Material)

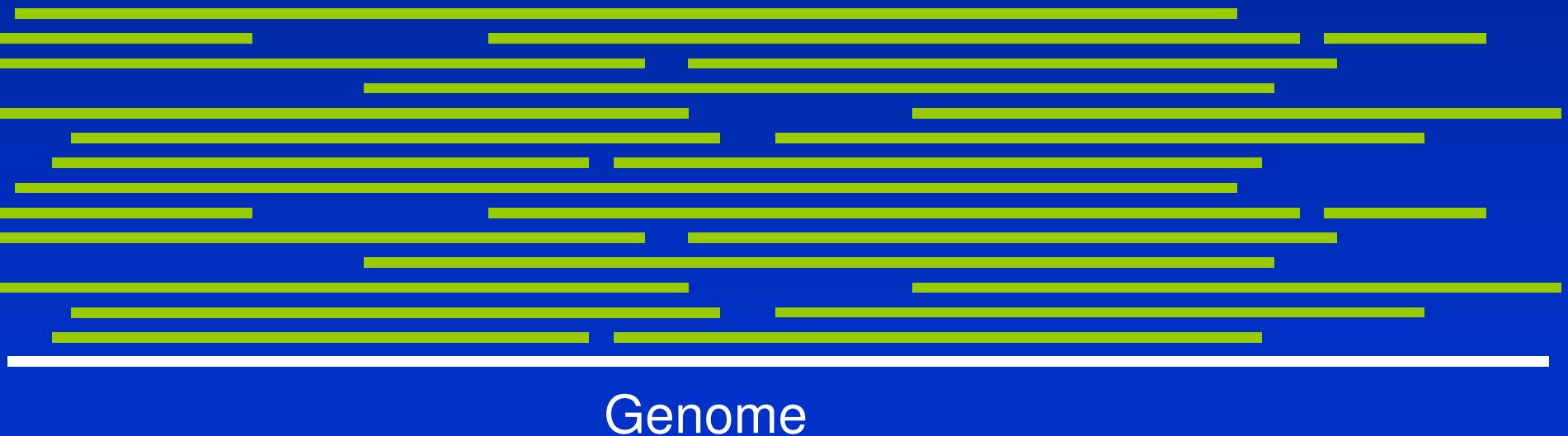


Heterozygosity and read depth

- SNP discovery
 - Missing some heterozygotes is not critical
 - Hopefully picked up in other animals
 - Just do more animals to identify SNP
 - Animal genotype not used directly
- Genotype calling
 - Missing heterozygotes a problem because incorrect genotype in downstream analysis
 - Statistical methods can be used to correct incorrect genotype calls
 - *Use genotype probabilities, not best guess!*

Identification of variants

- Program SAMtools
- stacks aligned bam files of multiple animals
- Calls variants and calculates quality/confidence statistics for calls
- <http://samtools.sourceforge.net/mpileup.shtml>



Variants in sequence

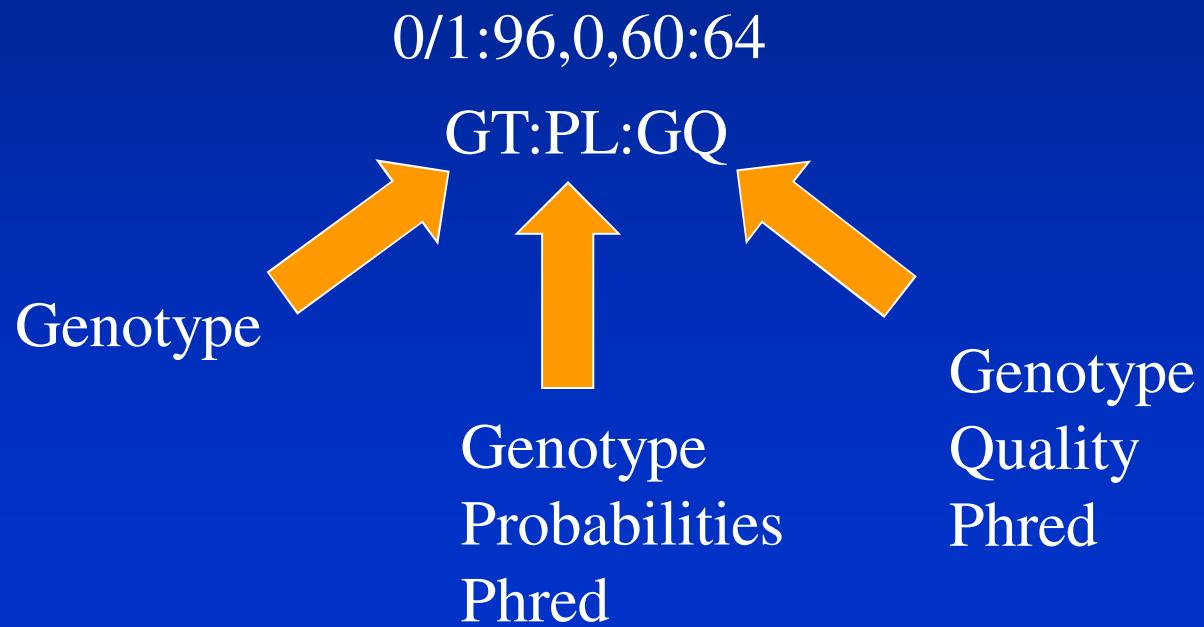
- SNP
- INDEL
 - INsertions and DEletions of DNA sections
- Copy number variants (CNV)
 - Repeated sections of DNA of various lengths
- Most studies to date have concentrated on SNP

VCF file

VCF file content							
<pre>##fileformat=VCFv4.1 ##samtoolsVersion=0.1.18 (1982:295) ##INFO=<ID=DP,Number=1,Type=Integer,Description="Raw read depth"> ##INFO=<ID=DP4,Number=4,Type=Integer,Description="# high-quality ref-forward bases, ref-reverse, alt-forward and alt-reverse bases"> ##INFO=<ID=MQ,Number=1,Type=Integer,Description="Root-mean-square mapping quality of covering reads"> ##INFO=<ID=FQ,Number=1,Type=Float,Description="Phred probability of all samples being the same"> ##INFO=<ID=AF1,Number=1,Type=Float,Description="Max-likelihood estimate of the first ALT allele frequency (assuming HWE)"> ##INFO=<ID=AC1,Number=1,Type=Float,Description="Max-likelihood estimate of the first ALT allele count (no HWE assumption)"> ##INFO=<ID=G3,Number=3,Type=Float,Description="ML estimate of genotype frequencies"> ##INFO=<ID=HWE,Number=1,Type=Float,Description="Chi^2 based HWE test P-value based on G3"> ##INFO=<ID=CLR,Number=1,Type=Integer,Description="Log ratio of genotype likelihoods with and without the constraint"> ##INFO=<ID=UGT,Number=1,Type=String,Description="The most probable unconstrained genotype configuration in the trio"> ##INFO=<ID=CGT,Number=1,Type=String,Description="The most probable constrained genotype configuration in the trio"> ##INFO=<ID=PV4,Number=4,Type=Float,Description="P-values for strand bias, baseQ bias, mapQ bias and tail distance bias"> ##INFO=<ID=INDEL,Number=0,Type=Flag,Description="Indicates that the variant is an INDEL."> ##INFO=<ID=PC2,Number=2,Type=Integer,Description="Phred probability of the nonRef allele frequency in group1 samples being larger (,smaller) than in group2."> ##INFO=<ID=PCHI2,Number=1,Type=Float,Description="Posterior weighted chi^2 P-value for testing the association between group1 and group2 samples."> ##INFO=<ID=QCHI2,Number=1,Type=Integer,Description="Phred scaled PCHI2.> ##INFO=<ID=PR,Number=1,Type=Integer,Description="# permutations yielding a smaller PCHI2.> ##INFO=<ID=VDB,Number=1,Type=Float,Description="Variant Distance Bias"> ##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype"> ##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality"> ##FORMAT=<ID=GL,Number=3,Type=Float,Description="Likelihoods for RR,RA,AA genotypes (R=ref,A=alt)"> ##FORMAT=<ID=DP,Number=1,Type=Integer,Description="# high-quality bases"> ##FORMAT=<ID=SP,Number=1,Type=Integer,Description="Phred-scaled strand bias P-value"> ##FORMAT=<ID=PL,Number=G,Type=Integer,Description="List of Phred-scaled genotype likelihoods"></pre>							
#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO
			FORMAT	Individual1	Individual2		
Chr29	39484430	.	C	A	277	GT:PL:GQ	0/0:0,9,113:8
			DP=30;VDB=0.0178;AF1=0.4455;AC1=5;DP4=11,4,7,8;MQ=52;FQ=279;PV4=0.26,0.43,0.0066,0.11				
			0/1:96,0,60:64				
Chr29	39484455	.	TGGG	TGG	18.6		
			INDEL;DP=23;VDB=0.0316;AF1=0.2602;G3=0.75,1.412e-06,0.25;HWE=0.0458;AC1=2;DP4=12,3,4,2;MQ=54;FQ=19.8;PV4=0.6,1,6.4e-05,1			GT:PL:GQ	
			0/0:0,9,90:11	0/0:0,9,93:11			
Chr29	39484540	.	A	G	999		
			DP=44;VDB=0.0356;AF1=0.588;AC1=6;DP4=7,8,14,14;MQ=46;FQ=999;PV4=1,1,0.079,1			GT:PL:GQ	0/0:0,15,157:11
Chr29	39484790	.	T	A	408		
			DP=33;VDB=0.0381;AF1=0.6663;AC1=7;DP4=6,2,14,11;MQ=50;FQ=413;PV4=0.43,0.21,0.0055,0.31			GT:PL:GQ	0/0:0,9,85:5
Chr29	39484791	.	A	C	999		
			DP=33;VDB=0.0381;AF1=0.6663;AC1=7;DP4=6,2,13,11;MQ=50;FQ=999;PV4=0.42,1,0.0069,0.33			GT:PL:GQ	0/0:0,9,88:5
			0/1:0,0,0:3				

VCF file (genotype probabilities)

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	Individual1	Individual2
Chr29	39484430	.	C	A	277	.	LocusStats	GT:PL:GQ	0/0:0,9,113:8	0/1:96,0,60:64
Chr29	39484455	.	TGGG	TGG	18.6	.	INDEL	LocusStats	GT:PL:GQ	0/0:0,9,90:11
Chr29	39484540	.	A	G	999	.	LocusStats	GT:PL:GQ	0/0:0,15,157:11	0/1:101,0,81:83
Chr29	39484790	.	T	A	408	.	LocusStats	GT:PL:GQ	0/0:0,9,85:5	0/1:0,0,0:3
Chr29	39484791	.	A	C	999	.	LocusStats	GT:PL:GQ	0/0:0,9,88:5	0/1:0,0,0:3



VCF file (FORMAT - locus quality stats)

In field FORMAT

DP=30;	Read depth
VDB=0.0178;	Variant distance bias
AF1=0.4455;	Maximum likelihood estimate of 1 st alternative allele frequency
AC1=5;	ML estimate of 1 alternative allele count
DP4=11,4,7,8;	Number reads on: Ref-Forward, Ref-Reverse, Alt-Forward, Alt-Reverse
MQ=52;	Mapping quality
FQ=279;	Phred probability of all samples being the same
PV4=0.26,0.43,0.0066,0.11	P-values for strand bias, baseQ bias, mapQ bias and tail distance bias

Filtering of variants

Reasons for filters:

- Number of artefacts of the sequencing process that lead to falsely identified variants
- Little evidence for a variant
 - Quality scores low

Reasons against filters:

- Real variants may be lost
 - Low frequency SNP often have lower quality scores

Variant filters we use (vcf)

1. Read depth
 - Minimum read depth
 - Individual genotype calls will be low quality
 - Maximum read depth
 - Short reads of repetitive regions may be mapped to same locations causing massive read depth
2. Mapping quality
 - Low quality calls
3. Quality
 - Phred score
4. Multiple variants within 3bp window
 - Alignment errors and indels can cause shifts → call 2 SNP close together instead of 1
 - Remove SNP close to indels

Phred quality scores (Q)

- Related to base-calling error probabilities.
Expressed in a range from 0 to 999 in our data.
- Probabilities are calculated by the following formula: $P = 10^{\frac{-Q}{10}}$
- e.g. Phred of 30 = error rate of 0.001
- Phred of 20 = error rate of 0.01
- Result is probability of each genotype at each variant eg. AA=0.95 AT=0.05 TT=0.00
- Use these in BEAGLE!

Imputation of full sequence data

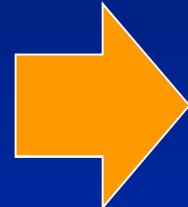
Create BAM files

1. Filter reads on quality score, trim ends
2. Remove PCR duplicates
3. Align with BWA



Variant calling

SamTools mPileup
Vcf file -> filter
(number forward /reverse reads of each allele, read depth, quality, filter number of variants in 5bp window)



Beagle Phasing in Reference

Input genotype probs from Phred scores
QC with 800K

Differences between SNP chip and sequence

- SNP chip
 - Sample of SNP
 - Higher minor allele frequency
 - Limited linkage disequilibrium depending on number of SNP
- Sequence
 - Contains most variants
 - SNP, indels, CNVs, etc
 - Allele frequency matches underlying causative variant frequency
 - Causative variants included
 - High linkage disequilibrium between variants

Using sequence data in genomic selection and GWAS

- Motivation
- Characteristics of sequence data
- Which individuals to sequence?
- Imputation of full sequence
- Methods for genomic prediction with full sequence data
- Examples
 - GWAS in Rice, Cattle

Which individuals to sequence?

- Those which capture greatest genetic diversity?
- Select set of individuals which are likely to capture highest proportion of unique chromosome segments

Which individuals to sequence?

- Let total number of individuals in population be n , number of individuals that can be sequenced be m .
- **A** = average relationship matrix among n individuals, from pedigree

- An example A matrix.....

Pedigree

Animal	Sire	Dam
1	0	0
2	0	0
3	0	0
4	1	2
5	1	2
6	1	3

Animals 6 is a half sib of 4 and 5

	Animal 1	Animal 2	Animal 3	Animal 4	Animal 5	Animal 6
Animal 1	1					
Animal 2	0	1				
Animal 3	0	0	1			
Animal 4	0.5	0.5	0	1		
Animal 5	0.5	0.5	0	0.5	1	
Animal 6	0.5	0	0.5	0.25	0.25	1

Which individuals to sequence?

- Let total number of individuals in population be n , number of individuals that can be sequenced be m .
- \mathbf{A} = average relationship matrix among n individuals, from pedigree
- \mathbf{c} is a vector of size n , which for each animal has the average relationship to the population (eg. Sum up the elements of \mathbf{A} down the column for individual i , take mean)

Which individuals to sequence?

- If we choose a group of m animals for sequencing, how much of the diversity do they capture
- $\mathbf{p}_m = \mathbf{A}_m^{-1} \mathbf{c}_m$
 - Where \mathbf{A}_m is the sub matrix of \mathbf{A} for the m individuals, and \mathbf{c}_m is the elements of the \mathbf{c} vector for the m individuals
- Proportion of diversity = $\mathbf{p}_m' \mathbf{1}_n$

Which individuals to sequence?

- Example

Which individuals to sequence?

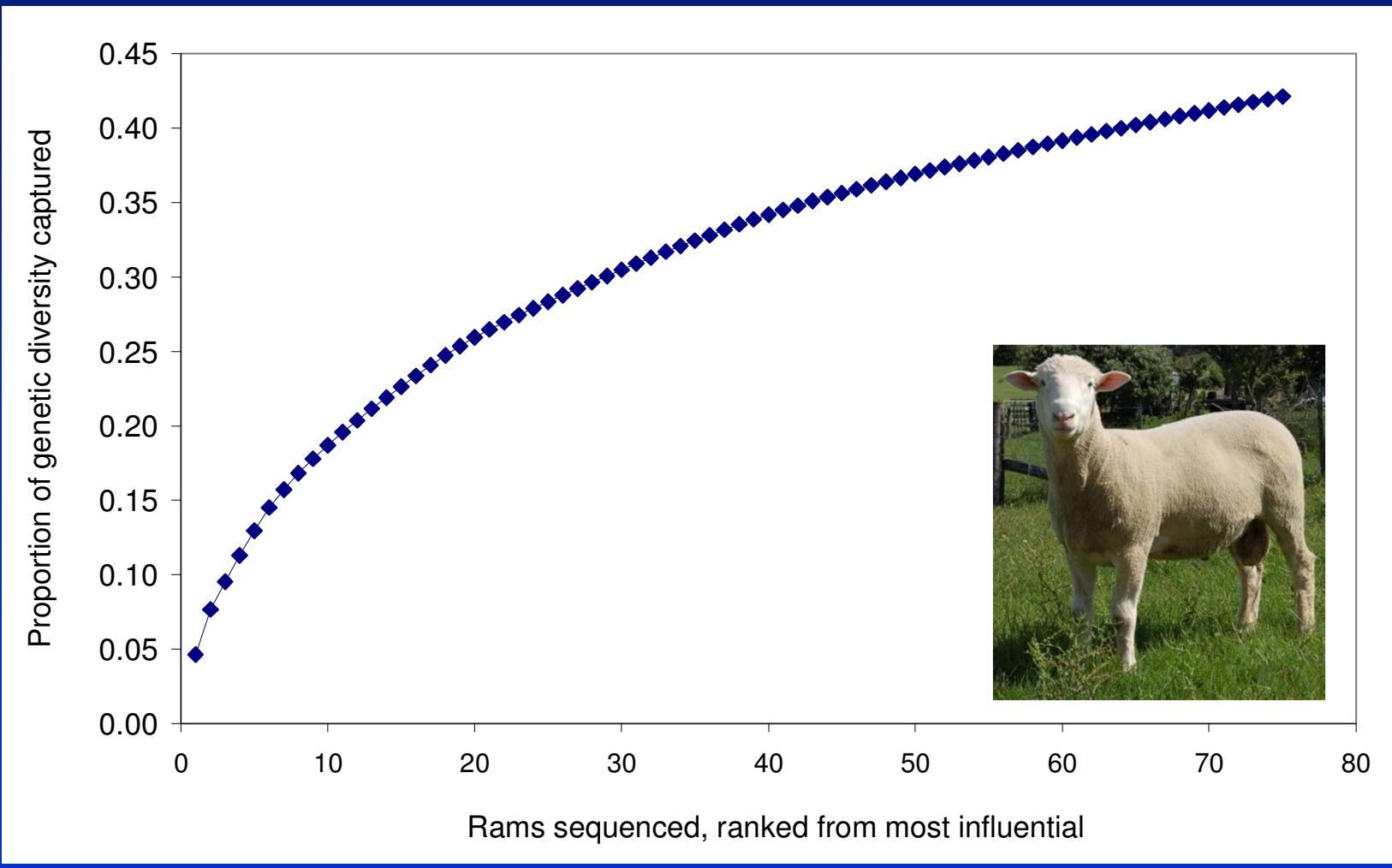
- Then choose set of individuals to sequence (m) which maximise $\mathbf{p}_m' \mathbf{1}_n$
- Step wise regression
 - Find single individual with largest p_i , set c_i to zero, next largest p_i , set c_i to zero.....
- Genetic algorithm

Which individuals to sequence?

- Then choose set of individuals to sequence (m) which maximise $\mathbf{p}_m' \mathbf{1} \mathbf{n}$
- Step wise regression
 - Find single individual with largest p_i , set c_i to zero, next largest p_i , set c_i to zero.....
- Genetic algorithm
- No **A**? Use **G**

Which individuals to sequence?

- Poll Dorset sheep



Using sequence data in genomic selection and GWAS

- Motivation
- Characteristics of sequence data
- Which individuals to sequence?
- Imputation of full sequence data
- Methods for genomic prediction with full sequence data
- Examples
 - GWAS in Rice, Cattle

Imputation of full sequence data

- Two groups of individuals
 - Sequenced individuals: reference population
 - Individuals genotyped on SNP array: target individuals

Imputation of full sequence data

- Steps:
 - Step 1. Find polymorphisms in sequence data
 - Step 2. Genotype all sequenced animals for polymorphisms (SNP, Indels)
 - Step 3. Phase genotypes (eg Beagle) in sequenced individuals, create reference file
 - Step 4. Impute all polymorphisms into individuals genotyped with SNP array

Imputation of full sequence data

Create BAM files

1. Filter reads on quality score, trim ends
2. Remove PCR duplicates
3. Align with BWA

BAM

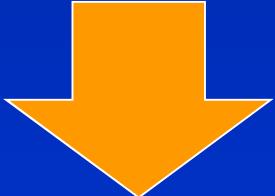
Variant calling

SamTools mPileup
Vcf file -> filter
(number forward /reverse reads of each allele, read depth, quality, filter number of variants in 5bp window)

Beagle Phasing in Reference

Input genotype probs from Phred scores
QC with 800K

Reference file for imputation



Analysis

- Genome wide association
Genomic selection

Genotype probabilities

Beagle Imputation in Target

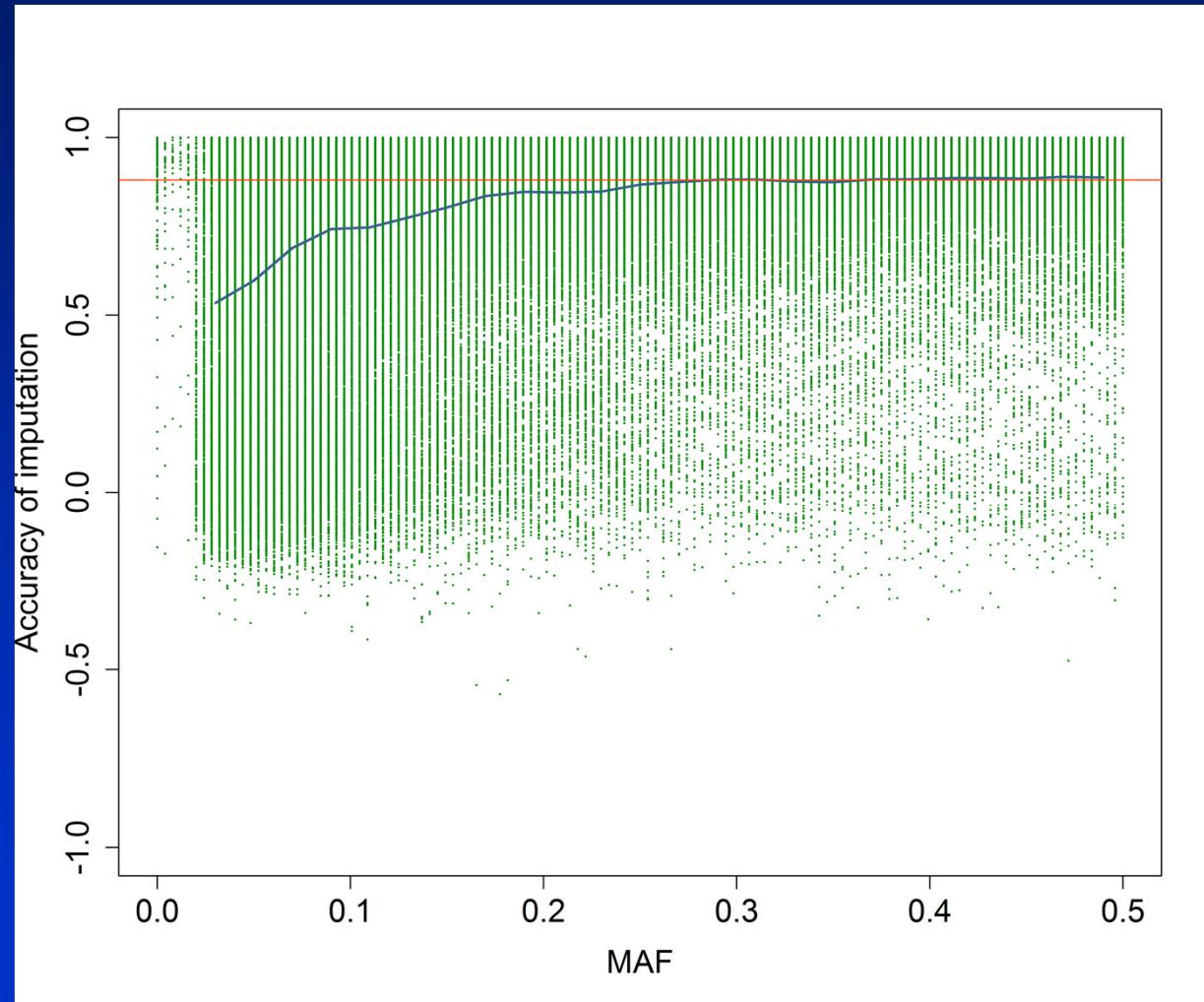
SNP array data in target population

Imputation of full sequence data

- How accurate?

Imputation

- 133 Holsteins in reference



Van Binsbergen et al., 2013

Using sequence data in genomic selection and GWAS

- Motivation
- Characteristics of sequence data
- Which individuals to sequence?
- Imputation of full sequence data
- Methods for genomic prediction with full sequence data
- Examples
 - GWAS in Rice, Cattle

Methods for genomic prediction with full sequence

- 14 million SNPs in Holstein Friesian cattle?
- Which method is most appropriate
- Priors
 - BLUP (GBLUP) -> all SNPs in LD with QTL, very small effects
 - BayesA -> some SNPs have moderate to large effects, rest very small
 - BayesB -> many SNPs have zero effect, some have small to moderate effect?

Methods for genomic prediction with full sequence

- Meuwissen and Goddard 2010
 - Simulated population with full sequence data, ~ 900 mutations chosen to be QTL
 - Used BLUP and BayesB to predict GEBV

The accuracy of the predictions of total genetic value (\pm SE) in the TEST1 data set when the training data contained $T = 200$ individuals and GWBLUP or BayesB is used to estimate the marker effects

Data	Causative SNPs			
	GWBLUP		BayesB	
	Excluded	Included	Excluded	Included
3 QTL	0.503 ± 0.011	0.508 ± 0.011	0.938 ± 0.013	0.973 ± 0.004
30 QTL	0.491 ± 0.016	0.493 ± 0.010	0.806 ± 0.023	0.826 ± 0.019

Methods for genomic prediction with full sequence

- Meuwissen and Goddard 2010
 - Simulated population with full sequence data, ~ 900 mutations chosen as QTL
 - Used BLUP and BayesB to predict GEBV
 - Large advantage of BayesB over BLUP
 - Prior matches their simulated data -> only 900 QTL amongst millions of SNP
 - 3% advantage of having mutation in data
 - Real data??

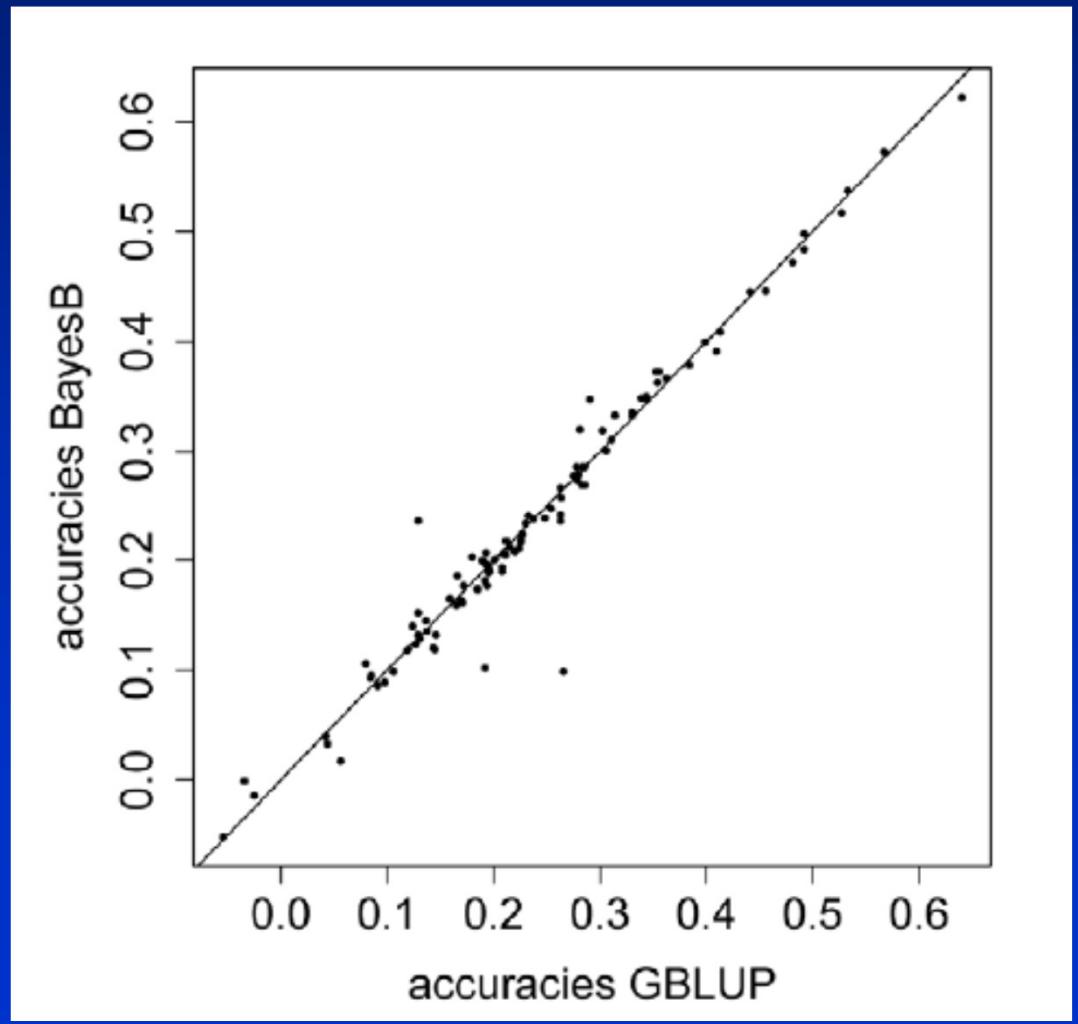
Methods for genomic prediction with full sequence

- Meuwissen and Goddard 2010
 - Better persistence of accuracy over generations

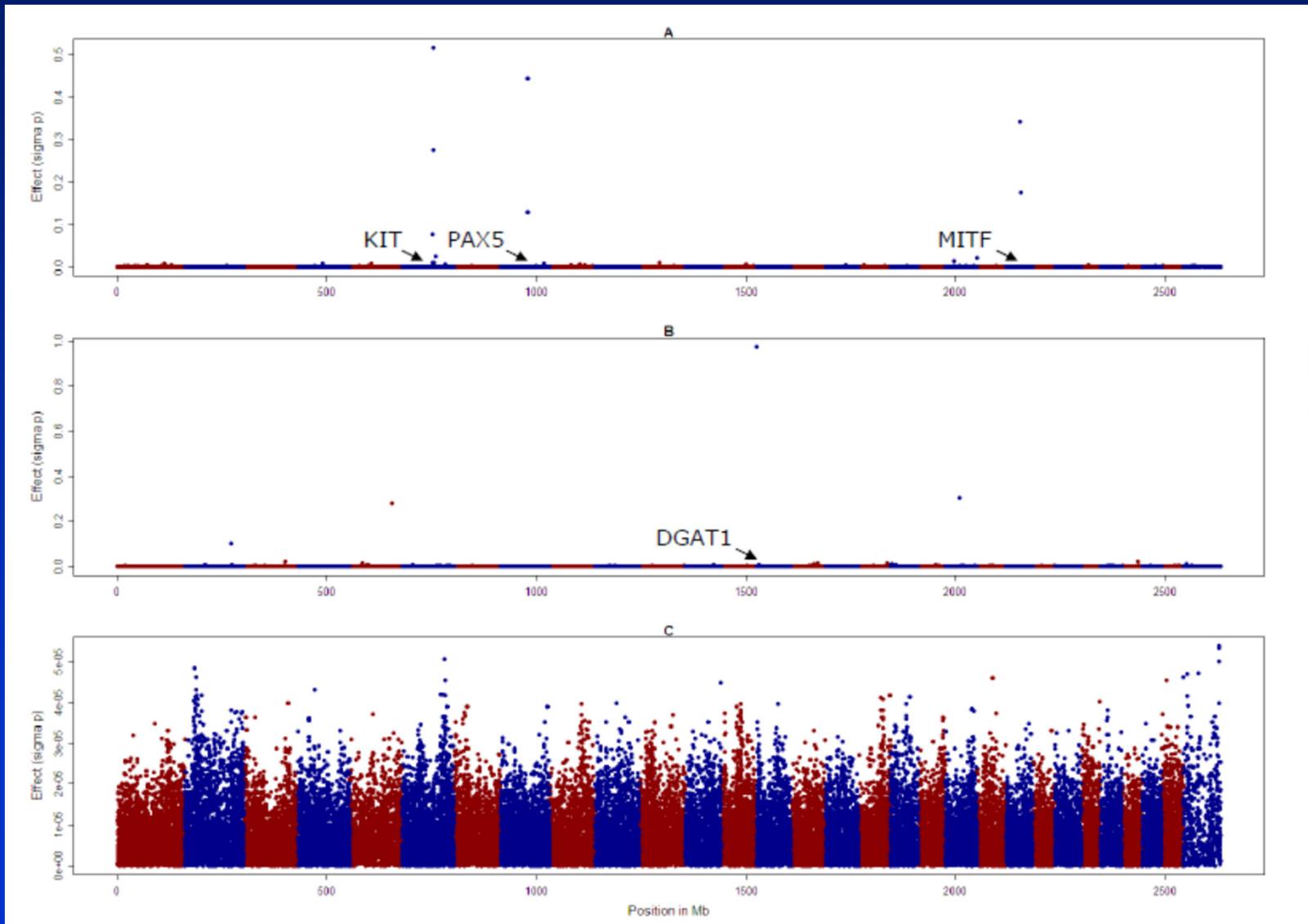
Causal SNPs	TEST1: $T = 200, L = 1:$ 30 QTL	TEST2: $T = 200, L = 1:$ 30 QTL
Excluded	0.806 ± 0.023	0.806 ± 0.022
Included	0.826 ± 0.019	0.824 ± 0.019

Methods for genomic prediction with full sequence

- Ober et al (2012) PLoS Genetics 8(5): e1002685
- Sample size
 - 157 fly lines
- No difference
 - BLUP vs BayesB



Genomic selection methods for GWAS?



Using sequence data in genomic selection and GWAS

- Motivation
- Characteristics of sequence data
- Which individuals to sequence?
- Imputation of full sequence data
- Methods for genomic prediction with full sequence data
- Examples
 - GWAS in Rice, Cattle

GWAS with sequence

ARTICLES

nature
genetics

Genome-wide association studies of 14 agronomic traits in rice landraces

Xuehui Huang^{1,2,10}, Xinghua Wei^{3,10}, Tao Sang^{4,10}, Qiang Zhao^{1,2,10}, Qi Feng^{1,10}, Yan Zhao¹, Canyang Li¹, Chuanrang Zhu¹, Tingting Lu¹, Zhiwu Zhang⁵, Meng Li^{5,6}, Danlin Fan¹, Yunli Guo¹, Ahong Wang¹, Lu Wang¹, Liuwei Deng¹, Wenjun Li¹, Yiqi Lu¹, Qijun Weng¹, Kunyan Liu¹, Tao Huang¹, Taoying Zhou¹, Yufeng Jing¹, Wei Li¹, Zhang Lin¹, Edward S Buckler^{5,7}, Qian Qian³, Qi-Fa Zhang⁸, Jiayang Li⁹ & Bin Han^{1,2}

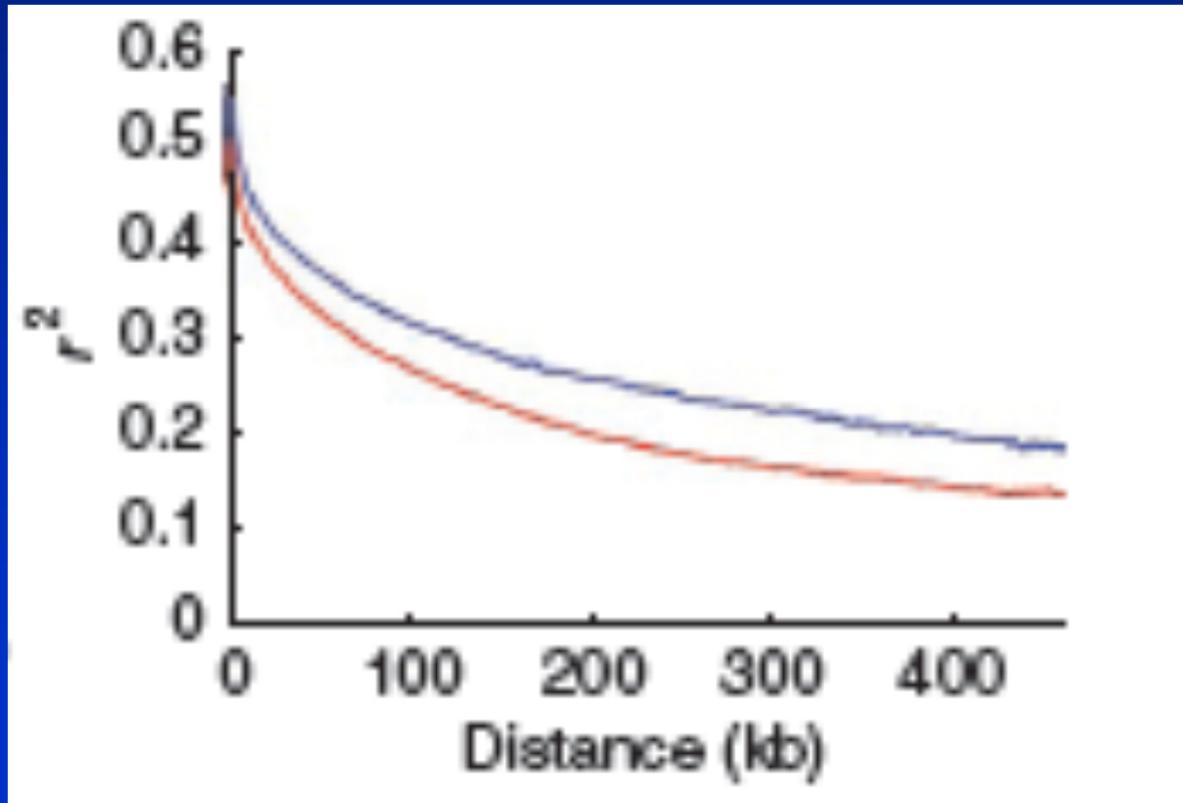
Uncovering the genetic basis of agronomic traits in crop landraces that have adapted to various agro-climatic conditions is important to world food security. Here we have identified ~3.6 million SNPs by sequencing 517 rice landraces and constructed a high-density haplotype map of the rice genome using a novel data-imputation method. We performed genome-wide association studies (GWAS) for 14 agronomic traits in the population of *Oryza sativa indica* subspecies. The loci identified through GWAS explained ~36% of the phenotypic variance, on average. The peak signals at six loci were tied closely to previously identified genes. This study provides a fundamental resource for rice genetics research and breeding, and demonstrates that an approach integrating second-generation genome sequencing and GWAS can be used as a powerful complementary strategy to classical biparental cross-mapping for dissecting complex traits in rice.

GWAS with sequence

- Huang et al. (2010)
 - Sequenced 517 rice landraces (inbred lines!) at 1x coverage
 - Represent ~ 82% of diversity in worlds rice cultivars
 - Called SNP in sequence pileups
 - 3.6 million SNP
 - With 1x coverage, could only call genotypes at ~ 20% of SNP
 - Therefore use imputation to fill in missing genotype
 - Example

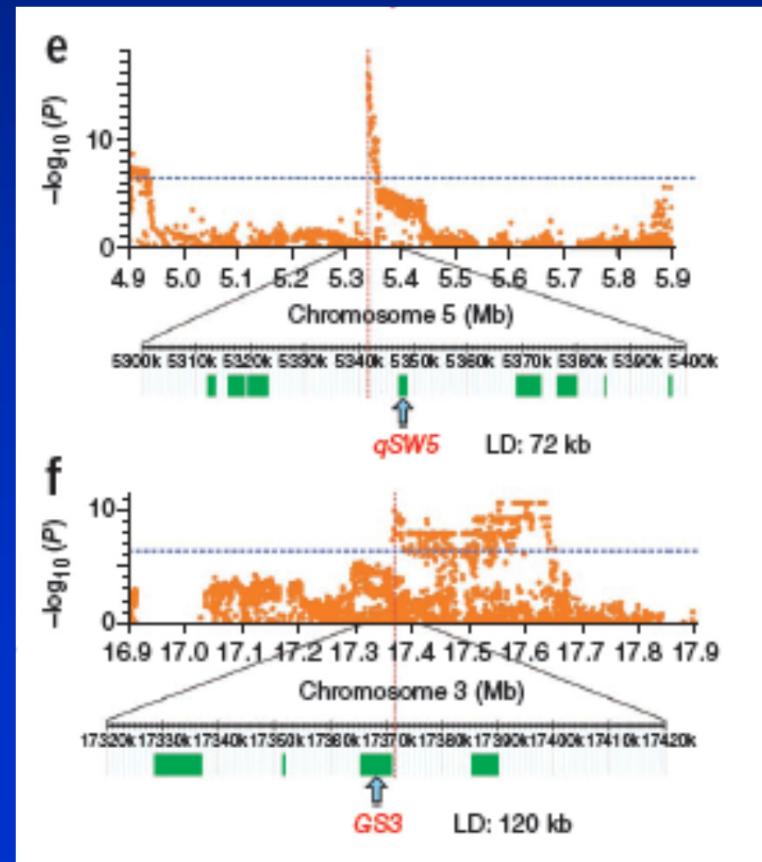
GWAS with sequence

- Huang et al. (2010)
 - Extent of LD
 - Red indica, blue japonica



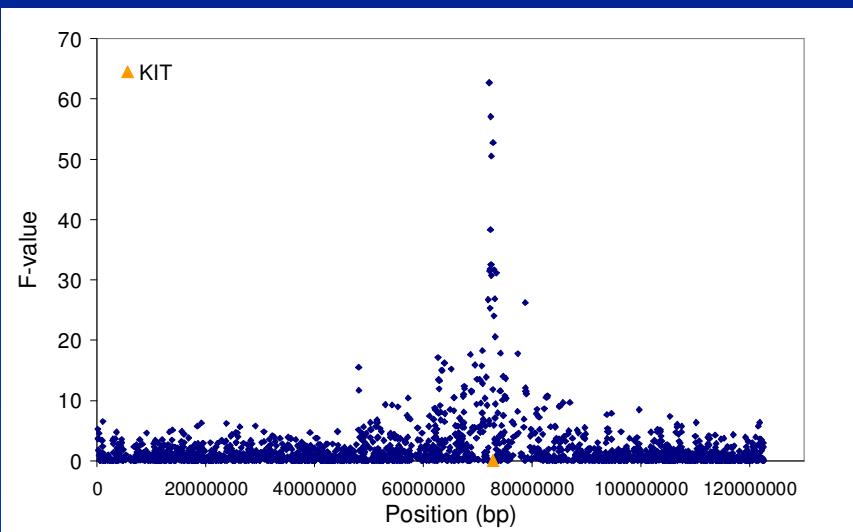
GWAS with sequence

- Huang et al. (2010)
 - Now have 517 lines with 3.6 million SNP genotyped
 - Well characterised phenotypes for 14 agronomic traits
 - Grain size, flowering date, etc
- Perform GWAS!
- Confirmed known mutations
- Many new mutations



GWAS with sequence

- KIT example
 - Earlier genome wide association study for proportion of black in Holsteins found association with SNP in KIT locus



- Can we impute sequence in this region and re-run association study?

GWAS with sequence

	Average fold coverage	Filtered SNPs	Concordance with 800K
PICKARD-ACRES VIC KAI	10.4	3,061,950	
GLENAGFTON ENHANCER	10.9	2,934,805	99.9%
BUSHLEA WAVES FABULON	11.3	4,249,998	97.4%
HANOVERHILL STARBUCK	12.5	3,237,681	97.9%
BIS-MAY S-E-L MOUNTAIN ET	12.6	3,009,463	98.5%
SHOREMAR PERFECT STAR	13.6	2,985,205	
ROYBROOK STARLITE	14.9	3,421,859	97.6%
TOPSPEED H POTTER	15.0	3,839,627	
LOCHAVON RAMESES	16.2	3,986,520	
BRAE DALE GOLDWYN	17.2	3,559,227	97.9%
CARENDA GRAVITY	17.8	4,331,849	96.8%
ONKAVALE GRIFFLAND MIDAS	22.5	3,742,799	

Imputation of full sequence data

Create BAM files

1. Filter reads on quality score, trim ends
2. Remove PCR duplicates
3. Align with BWA

BAM

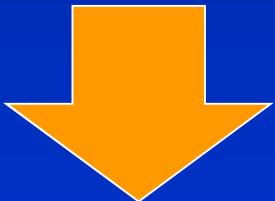
Variant calling

SamTools mPileup
Vcf file -> filter
(number forward /reverse reads of each allele, read depth, quality, filter number of variants in 5bp window)

Beagle Phasing in Reference

Input genotype probs from Phred scores
QC with 800K

Reference file for imputation



Analysis

- Genome wide association
Genomic selection

Genotype probabilities

Beagle Imputation in Target

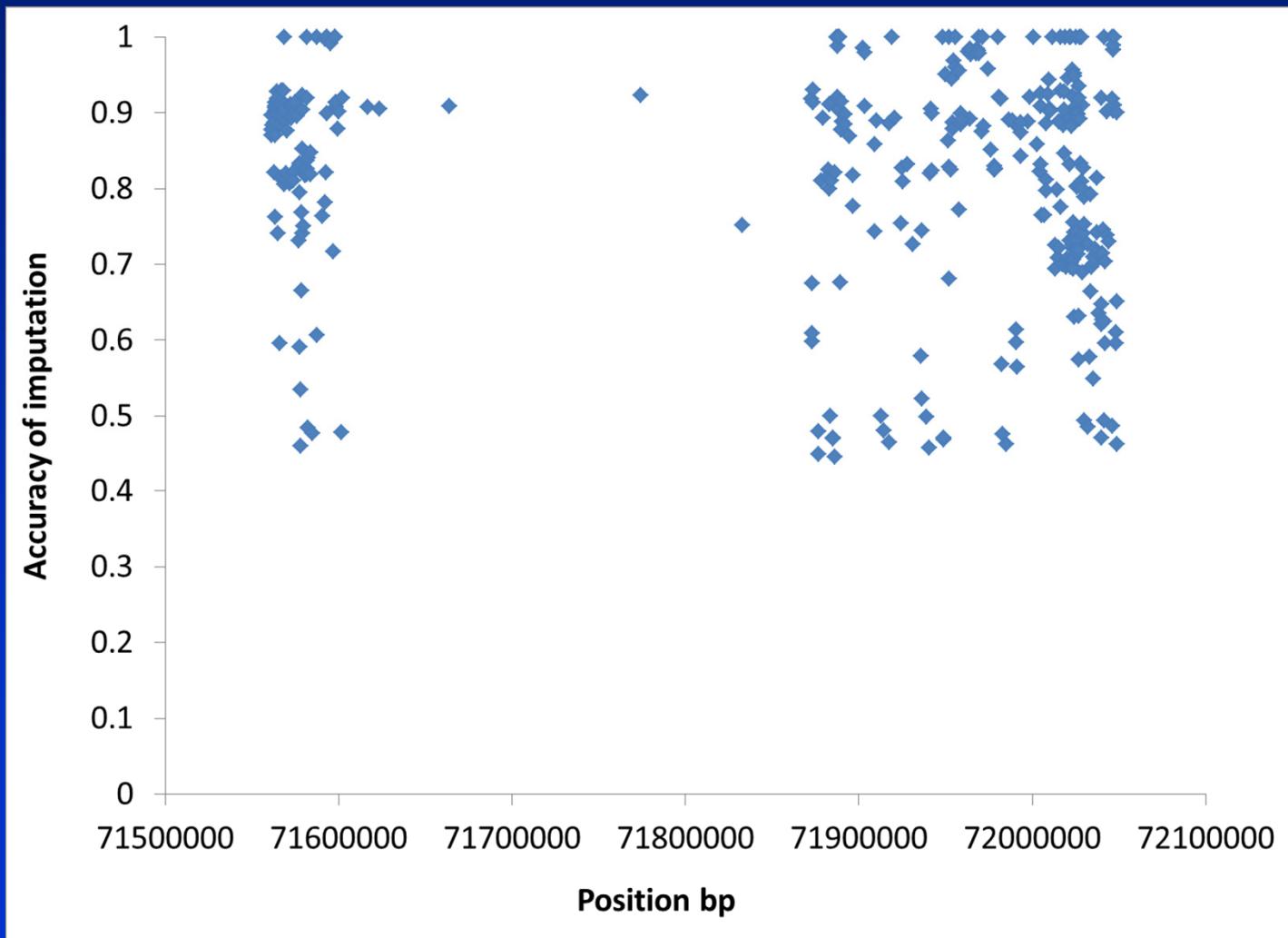
SNP array data in target population

GWAS with sequence

- KIT example
 - In sequenced bulls, compile list of SNPs/Indels in KIT region (352/20)
 - Call genotypes for the 372 variants in the 12 bulls
 - Use this as reference file for imputing the 372 variants in 697 bulls with % black phenotype (from 800K) data
 - Run association study on the 372 variants imputed in 697 bulls

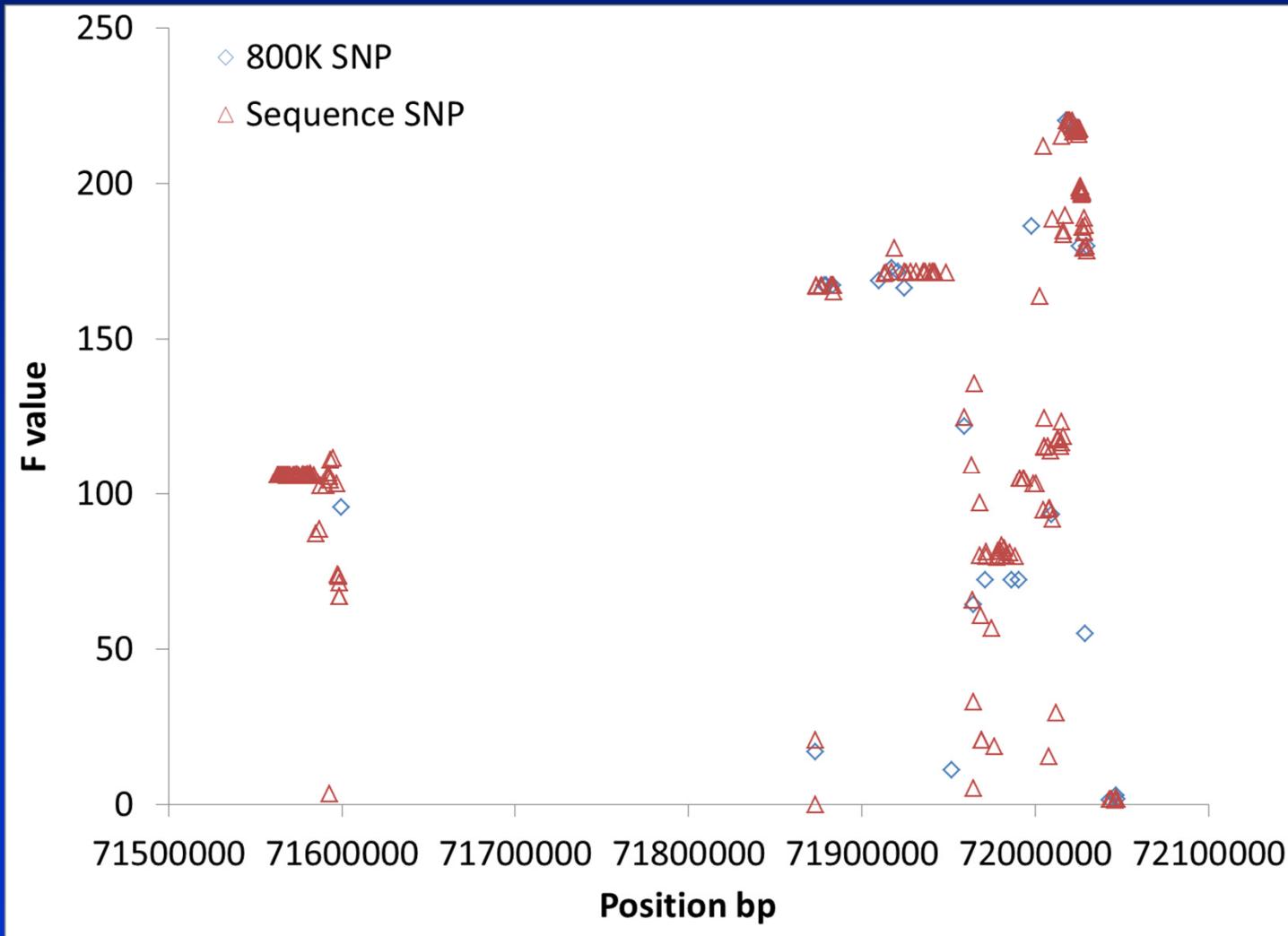
GWAS with sequence

- KIT example



GWAS with sequence

- KIT example



1000 Bull Genomes Project

- We will all need “reference” population of many sequenced bulls to impute from
 - SNP, indel and CNV genotypes
 - The more bulls the better!
- Combine sequence files (BAM) from around the world
- Run SNP/indel/CNV calling software several times a year
- Run 2.0 identified 28.3 million variants
- Contributors can download SNP/indel/CNV genotype file on all bulls to use for imputation anytime
- Partners welcome!

GWAS with sequence

- An alternative approach to GWAS?
 - For a target QTL region, sequence bulls of known QTL genotype (eg QQ,Qq,qq)
 - Have converted complex trait into a Mendelian trait
 - Far fewer individuals required for same power
 - Requires knowledge from linkage studies/previous GWAS!
 - Which method is more successful?

Quality of reference genomes?

- Cattle
 - Bovine build 4.2
 - Annotated
 - But many genes no assigned function
 - No Y chromosome yet, X is messy
 - ~ 9.5 million putative SNP in dbSNP
- Map of copy number variation?
- Kijas et al. (2010) – 51 CNV detected, 82% spanned at least one gene
- Hou et al. (2011) – 682 CNV from SNP array intensity data



Conclusions

- Potential of whole genome sequence data
 - Enable genome wide association study -> straight to causative mutation
 - Genomic selection
 - No longer have to rely on LD, causative mutation actually in data set, Higher accuracy of prediction?, Better persistence of accuracy across generations

Conclusions

- Potential of whole genome sequence data
 - Enable genome wide association study -> straight to causative mutation
 - Genomic selection
 - No longer have to rely on LD, causative mutation actually in data set, Higher accuracy of prediction?, Better persistence of accuracy across generations
- Choose individuals to sequence based on genetic contribution to population?
- Imputation of target population genotyped with SNP arrays
 - Caution with low frequency alleles, relationship to reference
- Large collaborative projects required for bovine/plant communities?