

# Whole Genome Sequencing

## Low Pass Sequencing

Gonçalo Abecasis

# Previous Lecture

- Introduction to Whole Genome Sequencing
  - What will we learn from whole genome sequencing?
- Challenges with Read Mapping
- Interpreting Mismatches: Variant or Error
  - Single individual analyses require deep sequencing
  - Multi-individual analyses can use shallower data
- Information contained in paired reads

# Questions that Might Be Answered With Complete Sequence Data...

- What is the contribution of each identified locus to a trait?
  - Likely that multiple variants, common and rare, will contribute
- What is the mechanism? What happens when we knockout a gene?
  - Most often, the causal variant will not have been examined directly
  - Rare coding variants will provide important insights into mechanisms
- What is the contribution of structural variation to disease?
  - These are hard to interrogate using current genotyping arrays.
- Are there additional susceptibility loci to be found?
  - Only subset of functional elements include common variants ...
  - Rare variants are more numerous and thus will point to additional loci

# Shotgun Sequence Data



TAGCTGATAGCTAG**A**TAGCTGATGAGCCCGAT

ATAGCTAG**A**TAGCTGATGAGCCCGATCGCTGCTAGCTC

ATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCC

AGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTG

GCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGA

Sequence Reads

5'-ACTGGTCGATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

$P(\text{reads} | A/A, \text{read mapped}) = 0.00000098$

$P(\text{reads} | A/C, \text{read mapped}) = 0.03125$

$P(\text{reads} | C/C, \text{read mapped}) = 0.000097$

Combine these likelihoods with a prior incorporating information from other individuals and flanking sites to assign a genotype.

# From Sequence to Genotype: Individual Based Prior



TAGCTGATAGCTAG**A**TAGCTGATGAGCCCGAT

ATAGCTAG**A**TAGCTGATGAGCCCGATCGCTGCTAGCTC

ATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCC

AGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTG

GCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGA

Sequence Reads

5'-ACTGGTCGATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

$P(\text{reads} | A/A) = 0.00000098$      $\text{Prior}(A/A) = 0.00034$      $\text{Posterior}(A/A) = <.001$

$P(\text{reads} | A/C) = 0.03125$      $\text{Prior}(A/C) = 0.00066$      $\text{Posterior}(A/C) = 0.175$

$P(\text{reads} | C/C) = 0.000097$      $\text{Prior}(C/C) = 0.99900$      $\text{Posterior}(C/C) = 0.825$

**Individual Based Prior:** Every site has 1/1000 probability of varying.

# From Sequence To Genotype: Population Based Prior



TAGCTGATAGCTAG**A**TAGCTGATGAGCCCGAT  
ATAGCTAG**A**TAGCTGATGAGCCCGATCGCTGCTAGCTC  
ATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCC  
AGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTG  
GCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGA

Sequence Reads

5'-ACTGGTCGATGCTAGCTGATAGCTAG**C**TAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

$P(\text{reads} | A/A) = 0.00000098$      $\text{Prior}(A/A) = 0.04$      $\text{Posterior}(A/A) = <.001$

$P(\text{reads} | A/C) = 0.03125$      $\text{Prior}(A/C) = 0.32$      $\text{Posterior}(A/C) = 0.999$

$P(\text{reads} | C/C) = 0.000097$      $\text{Prior}(C/C) = 0.64$      $\text{Posterior}(C/C) = <.001$

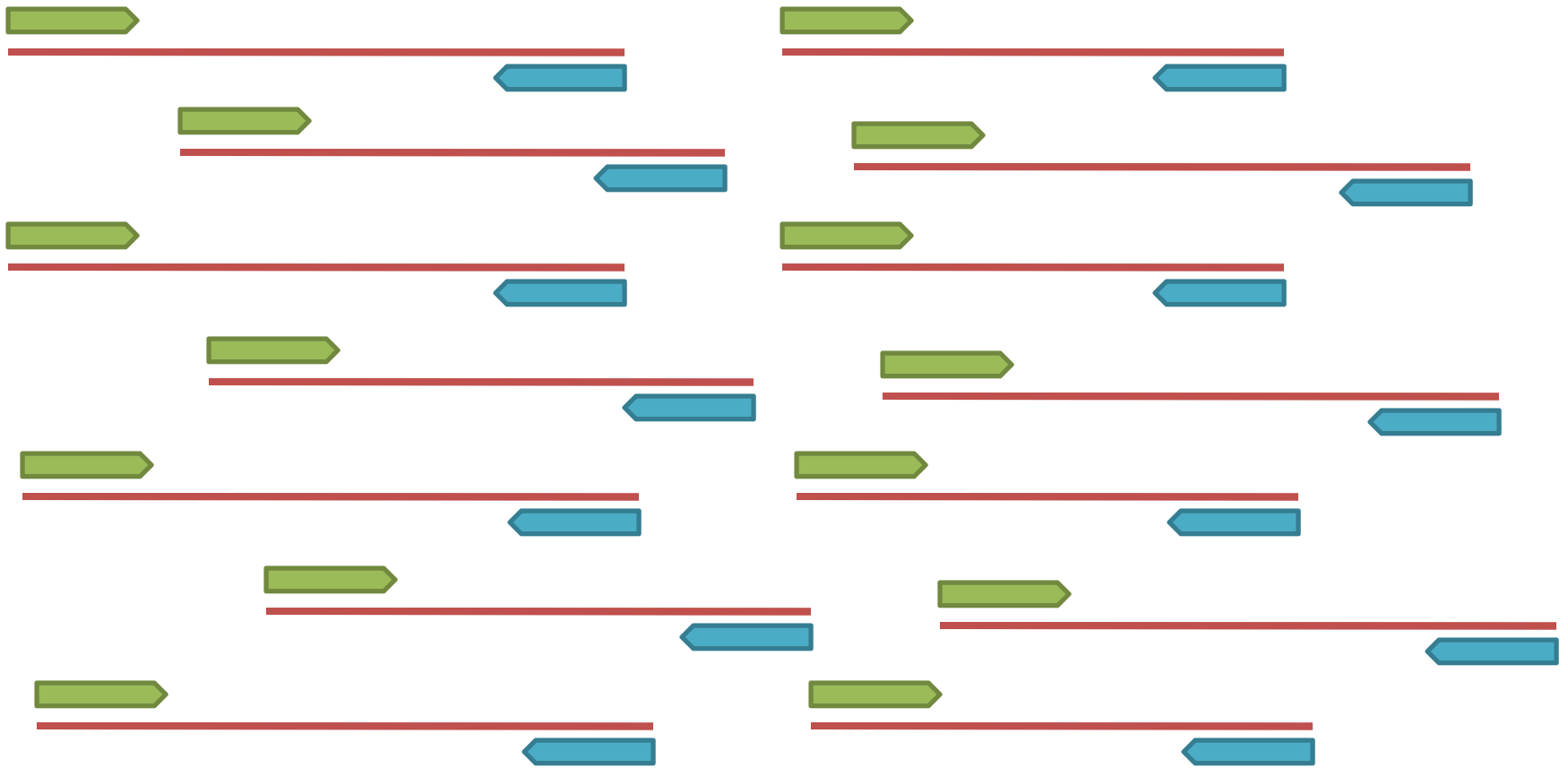
**Population Based Prior:** Use frequency information from examining others at the same site.

*In the example above, we estimated  $P(A) = 0.20$*

# Sequence Based Genotype Calls

- **Individual Based Prior**
  - Assumes all sites have an equal probability of showing polymorphism
  - Specifically, assumption is that about 1/1000 bases differ from reference
  - If reads were error free and sampling Poisson ...
  - ... 14x coverage would allow for 99.8% genotype accuracy
  - ... 30x coverage of the genome needed to allow for errors and clustering
- **Population Based Prior**
  - Uses frequency information obtained from examining other individuals
  - Calling very rare polymorphisms still requires 20-30x coverage of the genome
  - Calling common polymorphisms requires much less data
- **Haplotype Based Prior or Imputation Based Analysis**
  - Compares individuals with similar flanking haplotypes
  - Calling very rare polymorphisms still requires 20-30x coverage of the genome
  - Can make accurate genotype calls with 2-4x coverage of the genome
  - Accuracy improves as more individuals are sequenced

# Paired End Sequencing



Population of DNA fragments of known size (mean + stdev)  
Paired end sequences



# Paired End Sequencing

Paired Reads



Initial alignment to the reference genome



Paired end resolution



# Detecting Structural Variation

- Read depth
  - Regions where depth is different from expected
    - Expectation defined by comparing to rest of genome ...
    - ... or, even better, by comparing to other individuals
- Split reads
  - If reads are longer, it may be possible to find reads that span the structural variation
- Discrepant pairs
  - If we find pairs of reads that appear to map significantly closer or further apart than expected, could indicate an insertion or deletion
  - For this approach, “physical coverage” which is the sum of read length and insert size is key
- De Novo Assembly

# The Challenge

- Whole genome sequence data will greatly increase our understanding of complex traits
- Although a handful of genomes have been sequenced, this remains a relatively expensive enterprise
- Dissecting complex traits will require whole genome sequencing of 1,000s of individuals
- **How to sequence 1,000s of individuals cost-effectively?**

# Current Genome Scale Approaches

- Deep whole genome sequencing
  - Can only be applied to limited numbers of samples
  - Most complete ascertainment of variation
- Exome capture and targeted sequencing
  - Can be applied to moderate numbers of samples
  - SNPs and indels in the most interesting 1% of the genome
- Low coverage whole genome sequencing
  - Can be applied to moderate numbers of samples
  - Very complete ascertainment of shared variation
  - Less complete ascertainment of rare variants

# Current Genome Scale Approaches

- Deep whole genome sequencing
  - Can only be applied to limited numbers of samples
  - Most complete ascertainment of variation
- Exome sequencing
  - Can be applied to moderate numbers of samples
  - SNPs and indels in the most interesting 1% of the genome
- Low coverage whole genome sequencing
  - Can be applied to moderate numbers of samples
  - Very complete ascertainment of shared variation
  - Less complete ascertainment of rare variants

**Our Focus For Today**

# Recipe For Imputation With Shotgun Sequence Data

- Start with some plausible configuration for each individual
- Use Markov model to update one individual conditional on all others
- Repeat previous step many times
- Generate a consensus set of genotypes and haplotypes for each individual

# Silly Cartoon View of Shotgun Data

```
. G . G A . . T . C . T . T . . . T G .
C . A . . . C T C C C . . . C . . . . .
C C A . G . . C T . . . . . . . T G .
. . . . . . C T T T . C . . . . . . .
. . . . . T . . C . . A C C . . A T G .
. . . . . C . C C . G A C C . C A . G G
C G A . A . . . . . G . C . . T . T . .
. . . . . C . T . T . . . . . . A .
C G . . A . . C T . . . . . C T . G . .
C G A A . . T . . T . T . T . C T . . G C
. G A . A T C . . C . T . T T . . . G .
. . A . . . . . C C . A C . T C A T G .
. . A . G . . C . T T . . . T . T G . G C
C G A . . . T . . T . . . T T . T . . G C
. . . G A C . C . . . . . . . T G .
T . . . . T . . C . . . . . C C . . . .
. . . G A T C . C C . G . . C T T . . G C
. . . G A . T . T T . T . T T . T . . .
. G A G . . T . T . . G A . . T C G . . C
. . A A . . T . . . . . . . . . G .
```

# Cartoon View of Shotgun Data

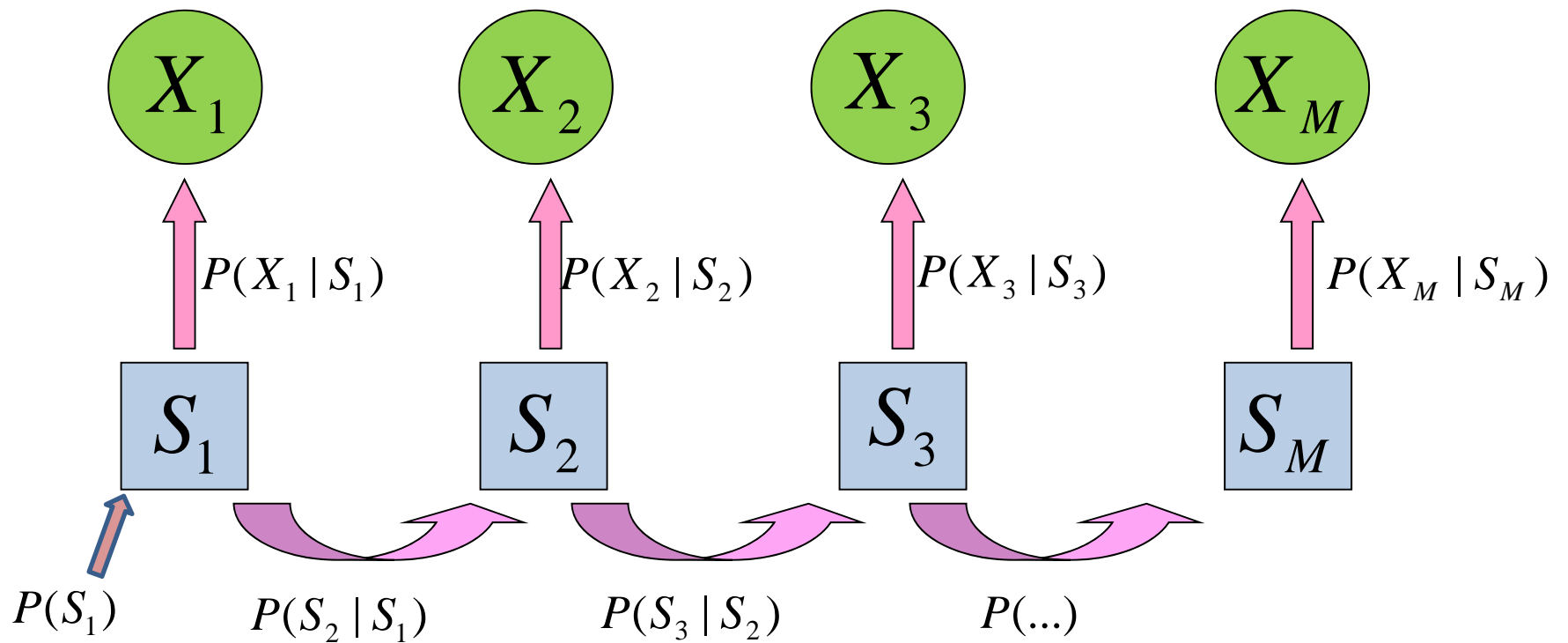
c	G	a	G	A	t	c	T	c	C	t	T	c	T	t	c	t	g	T	G	c	
C	g	A	g	a	t	C	T	C	C	C	g	a	c	C	t	c	a	t	g	g	
C	C	A	a	G	c	t	C	T	t	t	t	c	t	t	c	t	g	T	G	c	
c	g	a	a	g	c	t	C	T	T	T	t	C	t	t	c	t	g	t	g	c	
c	g	a	g	a	c	T	c	t	C	c	g	A	C	C	t	t	A	T	G	c	
t	g	g	g	a	t	C	t	C	C	c	G	A	C	C	t	C	A	t	G	G	
C	G	A	g	A	t	c	t	c	c	c	G	a	C	c	t	T	g	T	g	c	
c	g	a	g	a	c	t	C	t	T	t	T	c	t	t	t	t	g	t	A	c	
C	G	a	g	A	c	t	C	T	c	c	g	a	c	C	T	c	G	t	g	c	
C	G	A	A	g	c	T	c	t	T	t	T	c	T	t	C	T	g	t	G	C	
c	G	A	g	A	T	C	t	c	C	t	T	c	T	T	c	t	g	t	G	c	
c	g	A	g	a	t	c	t	c	C	C	g	A	C	c	T	C	A	T	G	g	
c	c	A	a	G	c	t	C	t	T	T	t	c	t	T	c	T	G	t	G	C	
C	G	A	a	g	c	T	c	t	T	t	t	c	T	T	c	T	g	t	G	C	
c	g	a	G	A	C	t	C	t	C	t	c	g	a	c	c	t	t	a	T	G	c
T	g	g	g	a	T	c	t	C	c	c	g	a	C	C	t	c	a	t	g	g	
c	g	a	G	A	T	C	t	C	C	c	G	a	c	C	T	T	g	t	G	C	
c	g	a	G	A	c	T	c	T	T	t	T	c	T	T	t	T	g	t	a	c	
c	G	A	G	a	c	T	c	T	c	c	G	A	c	c	T	C	G	t	g	C	
c	g	A	A	g	c	T	c	t	t	t	t	c	t	t	c	t	g	t	G	c	



# How Do We Update One Pair Of Haplotypes?

- Markov model is very similar to that used for analysis of genotype imputation analysis
- To carry out an update, select one individual
  - Let  $X_i$  be observed bases overlapping position  $i$  for individual
- Assume (temporarily) that current haplotype estimates for all other individuals are correct
- Model haplotypes for individual being updated as mosaic of the other available haplotypes
  - $S_i = (S_{i1}, S_{i2})$  denotes the pair of haplotypes being copied

# Markov Model



The final ingredient connects template states along the chromosome ...

# Likelihood

$$L = \sum_{S_1} \sum_{S_2} \dots \sum_{S_M} P(S_1) \prod_{i=2}^M P(S_i | S_{i-1}) \prod_{i=1}^M P(X_i | S_i)$$

- $P(S_1) = 1 / H^2$  where  $H$  is the number of template haplotypes
- $P(S_i | S_{i-1})$  depends on estimated population recombination rate
- $P(X_i | S_i)$  are the genotype likelihoods

# Simulation Results: Common Sites

- Detection and genotyping of Sites with MAF >5% (2116 simulated sites/Mb)
  - **Detected Polymorphic Sites: 2x coverage**
    - 100 people      2102 sites/Mb detected
    - 200 people      2115 sites/Mb detected
    - 400 people      2116 sites/Mb detected
  - **Error Rates at Detected Sites: 2x coverage**
    - 100 people      98.5% accurate, 90.6% at hets
    - 200 people      99.6% accurate, 99.4% at hets
    - 400 people      99.8% accurate, 99.7% at hets

# Simulation Results: Rarer Sites

- Detection and genotyping of Sites with MAF 1-2% (425 simulated sites/Mb)
  - **Detected Polymorphic Sites: 2x coverage**
    - 100 people      139 sites/Mb detected
    - 200 people      213 sites/Mb detected
    - 400 people      343 sites/Mb detected
  - **Error Rates at Detected Sites: 2x coverage**
    - 100 people      98.6% accurate, 92.9% at hets
    - 200 people      99.4% accurate, 95.0% at hets
    - 400 people      99.6% accurate, 95.9% at hets

**That's The Theory ...  
Show Me The Data!**

Results from 1000 Genomes Project

# Project Goals

- >95% of accessible genetic variants with a frequency of >1% in each of multiple continental regions
- Extend discovery effort to lower frequency variants in coding regions of the genome
- Define haplotype structure in the genome

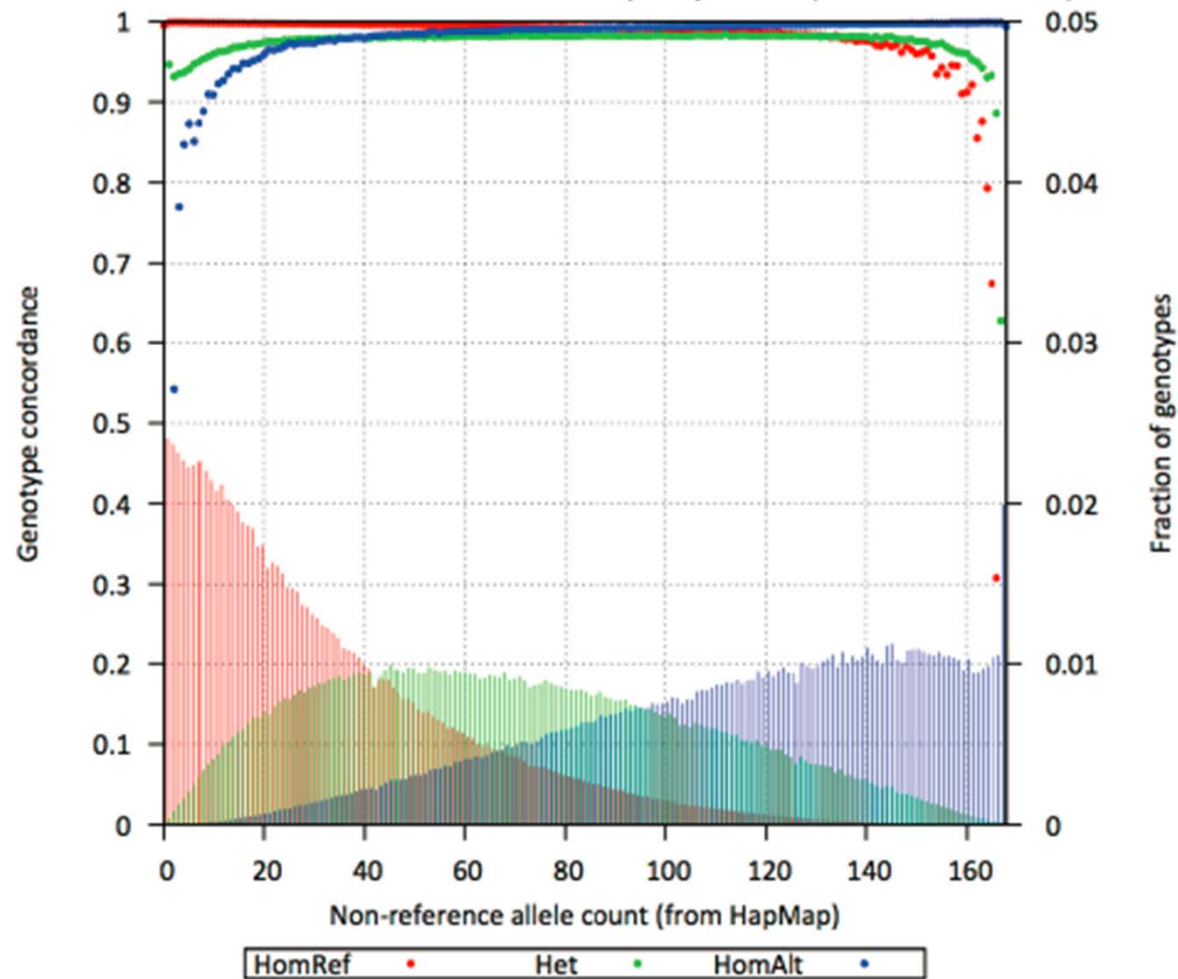
# 1000 Genomes Pilot Completed



- 2 deeply sequenced trios
- 179 whole genomes sequenced at low coverage
- 8,820 exons deeply sequenced in 697 individuals
- 15M SNPs, 1M indels, 20,000 structural variants



# Accuracy of Low Pass Genotypes



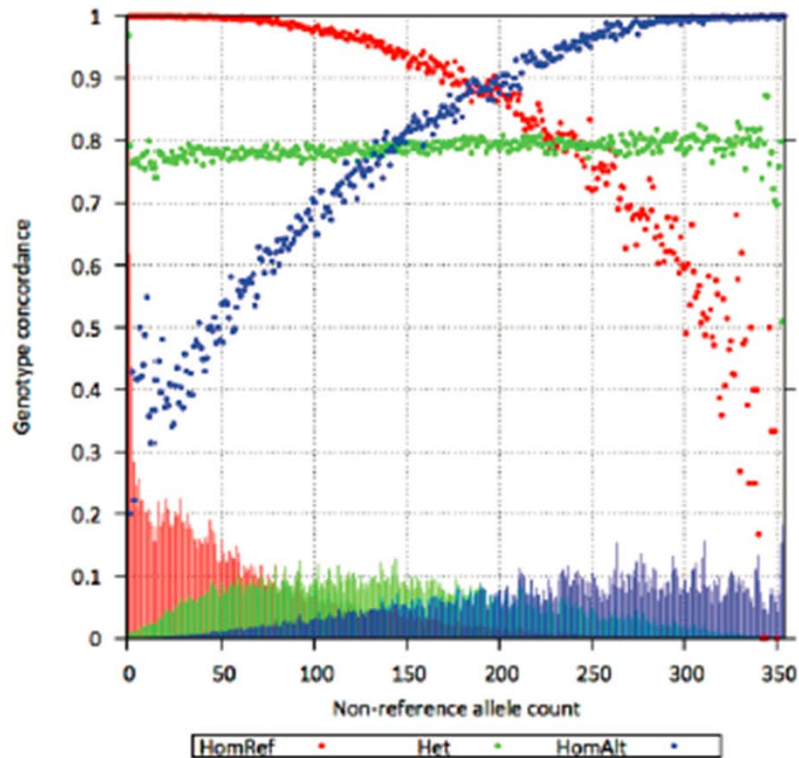
Genotype accuracy for rare genotypes is lowest, but definition of rare changes as more samples are sequenced.

Hyun Min Kang

# Does Haplotype Information Really Help?

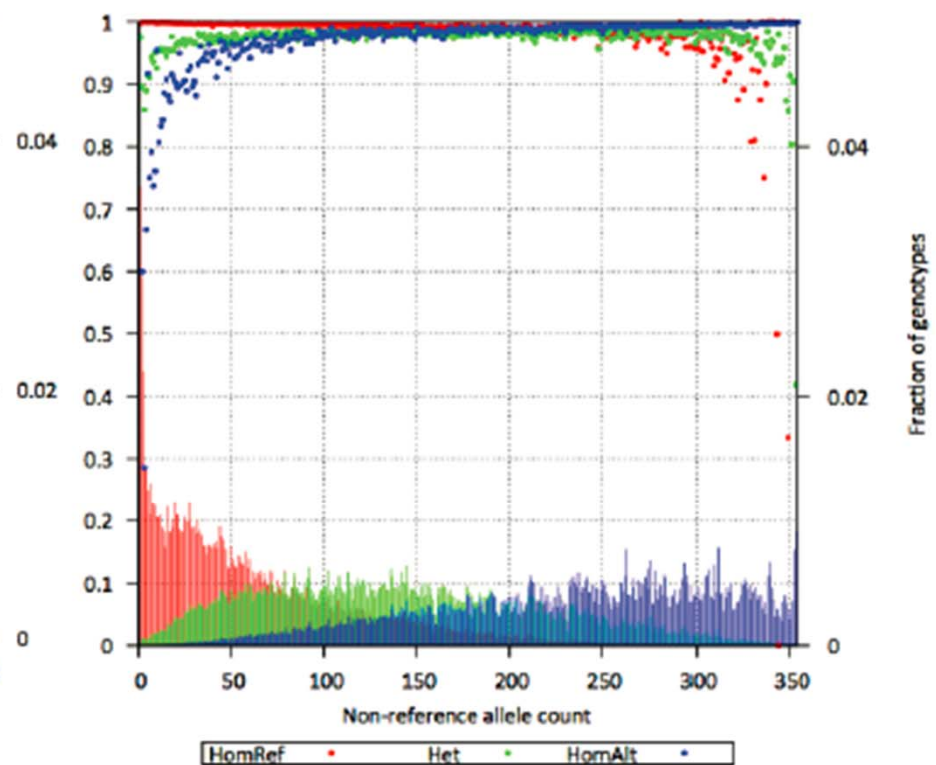
## Single Site Analysis

– 21.4% HET errors



## Haplotype Aware Analysis

– 2.0% HET errors



# As More Samples Are Sequenced, Low Pass Genotypes Improve

Analysis	#SNPs	dbSNP%	Missing HapMap %	Ts/Tv	Accuracy at Hets*
March 2010 Michigan/EUR 60	9,158,226	63.5	7.0	1.91	96.74
August 2010 Michigan/EUR 186	10,537,718	52.5	5.6	2.04	97.56
October 2010 Michigan/EUR 280	13,276,643	50.1	1.8	2.20	97.91**

Accuracy of Low Pass Genotypes Generated by 1000 Genomes Project,  
When Analysed At the University of Michigan

# Some Important Notes

- The Markov model we described is one of several possible models for analysis of low pass data
- Alternative models, based on E-M algorithms or local clustering of individuals into small groups exist
- Currently, the best possible genotypes produced by running multiple methods and generating a consensus across analysis their results.

# Implications for Whole Genome Sequencing Studies

- Suppose we could afford 2,000x data (6,000 GB)
- We could sequence 67 individuals at 30x

Sequencing of 67 individuals at 30x depth

Minor Allele Frequency	0.5 – 1.0%	1.0 – 2.0%	2.0 – 5.0%	>5%
Proportion of Detected Sites	59.3%	90.1%	96.9%	100.0%
Genotyping Accuracy	100.0%	100.0%	100.0%	100.0%
.... Heterozygous Sites Only	100.0%	100.0%	100.0%	100.0%
Correlation with Truth ( $r^2$ )	99.8%	99.9%	99.9%	100.0%
Effective Sample Size ( $n \cdot r^2$ )	67	67	67	67

# Implications for Whole Genome Sequencing Studies

- Suppose we could afford 2,000x data (6,000 GB)
- We could sequence 1000 individuals at 2x

Sequencing of 1000 individuals at 2x depth

Minor Allele Frequency	0.5 – 1.0%	1.0 – 2.0%	2.0 – 5.0%	>5%
Proportion of Detected Sites	79.6%	98.8%	100.0%	100.0%
Genotyping Accuracy	99.6%	99.5%	99.5%	99.8%
.... Heterozygous Sites Only	78.8%	89.5%	95.9%	99.8%
Correlation with Truth ( $r^2$ )	56.7%	76.1%	88.2%	97.8%
Effective Sample Size ( $n \cdot r^2$ )	567	761	882	978

# Summary for Today

- Analysis of Low Pass Sequence Data
  - Single sample analyses produce poor quality variants.
  - Single site analyses produce poor quality genotypes.
  - Multi-sample, multi-sample analyses can work quite well.
- Why low pass analyses are attractive for complex disease association studies.

# Recommended Reading

- The 1000 Genomes Project (2010) A map of human genome variation from population-scale sequencing. *Nature* **467**:1061-73
- Li Y, Willer CJ, Ding J, Scheet P and Abecasis GR (2010) MaCH: using sequence and genotype data to estimate haplotypes and unobserved genotypes. *Genet Epidemiol* **34**:816-834
- Le SQ and Durbin R (2010) SNP detection and genotyping from low-coverage sequencing data on multiple diploid samples. *Genome Research* (in press)