

# Variant Calling

Gonçalo Abecasis

Hyun Min Kang

# Calling Consensus Genotype - Details

- Each aligned read provides a small amount of evidence about the underlying genotype
  - Read may be consistent with a particular genotype ...
  - Read may be less consistent with other genotypes ...
  - A single read is never definitive
- This evidence is cumulated gradually, until we reach a point where the genotype can be called confidently
- Let's outline a simple approach ...

# 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

A/C

Predicted Genotype

# Shotgun Sequence Data

Sequence Reads

5'-ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

$P(\text{reads} | A/A, \text{read mapped}) = 1.0$

$P(\text{reads} | A/C, \text{read mapped}) = 1.0$

$P(\text{reads} | C/C, \text{read mapped}) = 1.0$

Possible Genotypes

# Shotgun Sequence Data

GCTAGCTGATAGCTAGCTAGCTAGCTGATGAGCCCGA

Sequence Reads

5'-ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

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

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

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

Possible Genotypes

# Shotgun Sequence Data

GCTAGCTGATAGCTAGCTAGCTAGCTGATGAGCCCGA

Sequence Reads

5'-ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

$P(\text{reads} | A/A, \text{read mapped}) = 0.01$

$P(\text{reads} | A/C, \text{read mapped}) = 0.50$

$P(\text{reads} | C/C, \text{read mapped}) = 0.99$

Possible Genotypes

# Shotgun Sequence Data



AGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTG  
GCTAGCTGATAGCTAGCTAGCTGATGAGCCCGA

Sequence Reads

5'-ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

$P(\text{reads} | A/A, \text{read mapped}) = 0.0001$

$P(\text{reads} | A/C, \text{read mapped}) = 0.25$

$P(\text{reads} | C/C, \text{read mapped}) = 0.98$

Possible Genotypes

# Shotgun Sequence Data

ATGCTAGCTGATAGCTAGCTAGCTAGCTGATGAGCC  
AGCTGATAGCTAGCTAGCTAGCTGATGAGCCCCGATCGCTG  
GCTAGCTGATAGCTAGCTAGCTAGCTGATGAGCCCCGA

Sequence Reads

5'-ACTGGTCGATGCTAGCTGATAGCTAGCTAGCTGATGAGCCCCGATCGCTGCTAGCTCGACG-3'

Reference Genome

$P(\text{reads} | A/A, \text{read mapped}) = 0.000001$

$P(\text{reads} | A/C, \text{read mapped}) = 0.125$

$P(\text{reads} | C/C, \text{read mapped}) = 0.97$

Possible Genotypes



# Shotgun Sequence Data



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

Sequence Reads

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

Reference Genome

$P(\text{reads} | A/A, \text{read mapped}) = 0.00000099$

$P(\text{reads} | A/C, \text{read mapped}) = 0.0625$

$P(\text{reads} | C/C, \text{read mapped}) = 0.0097$

Possible Genotypes

# 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

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

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

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

Possible Genotypes

# 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.

# 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

$$P(\text{Genotype}|\text{reads}) = \frac{P(\text{reads}|\text{Genotype})\text{Prior}(\text{Genotype})}{\sum_G P(\text{reads}|G)\text{Prior}(G)}$$

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

# Ingredients That Go Into Prior

- Most sites don't vary
  - P(non-reference base)  $\sim 0.001$
- When a site does vary, it is usually heterozygous
  - P(non-reference heterozygote)  $\sim 0.001 * 2/3$
  - P(non-reference homozygote)  $\sim 0.001 * 1/3$
- Mutation model
  - Transitions account for most variants (C $\leftrightarrow$ T or A $\leftrightarrow$ G)
  - Transversions account for minority of variants

# 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(reads   A/A) = 0.00000098	<b>Prior(A/A) = 0.00034</b>	Posterior(A/A) = <.001
P(reads   A/C) = 0.03125	<b>Prior(A/C) = 0.00066</b>	Posterior(A/C) = 0.175
P(reads   C/C) = 0.000097	<b>Prior(C/C) = 0.99900</b>	Posterior(C/C) = 0.825

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

# 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.

# 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



# 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

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

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

$P(\text{reads} | C/C) = 0.000097$     **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  $P(A) = 0.20$*

# 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

# Shotgun Sequence Data

## Haplotype 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

$P(\text{reads} | A/A) = 0.00000098$     **Prior(A/A) = 0.81**     $\text{Posterior}(A/A) = <.001$

$P(\text{reads} | A/C) = 0.03125$     **Prior(A/C) = 0.18**     $\text{Posterior}(A/C) = 0.999$

$P(\text{reads} | C/C) = 0.000097$     **Prior(C/C) = 0.01**     $\text{Posterior}(C/C) = <.001$

**Haplotype Based Prior:** Examine other chromosomes that are similar at locus of interest.

*In the example above, we estimated that 90% of similar chromosomes carry allele A.*

# Shotgun Sequence Data

## Haplotype 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

$P(\text{reads} | A/A) = 0.00000098$      $\text{Prior}(A/A) = 0.81$      $\text{Posterior}(A/A) = <.001$

$P(\text{reads} | A/C) = 0.03125$      $\text{Prior}(A/C) = 0.18$      $\text{Posterior}(A/C) = 0.999$

$P(\text{reads} | C/C) = 0.000097$      $\text{Prior}(C/C) = 0.01$      $\text{Posterior}(C/C) = <.001$

**Haplotype Based Prior:** Examine other chromosomes that are similar at locus of interest.

*In the example above, we estimated that 90% of similar chromosomes carry allele A.*

# 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

# Why Like Low Pass Sequencing? Whole Genome Sequencing Studies

	.5 – 1%	1 – 2%	2-5%
<b>400 Deep Genomes (30x)</b>			
Discovery Rate	100%	100%	100%
Het. Accuracy	100%	100%	100%
Effective N	400	400	400
<b>3000 Shallow Genomes (4x)</b>			
Discovery Rate	100%	100%	100%
Het. Accuracy	90.4%	97.3%	98.8%
Effective N	2406	2758	2873

# Variant Filtering

- Typically, an initial set of variant calls might include many artifacts
- Artifacts can be identified through a series of unusual features:
  - Clusters of SNPs
  - Unusual flanking sequences, particularly homopolymer runs
  - Unusual balance of reads with forward / reverse orientation
  - Unusual fraction of reads with reference base in heterozygotes
- Typically, one compares features like allele frequency or transition / transversion ratio for different sets of potential filters
  - For example, if we find that the transition/transversion ratio is unusually low when there is another SNP within the flanking 5-bp, we might decide to filter those.



# Example of Raw and Filtered Callsets

CATEGORY	# SNPs	#dbSNP (b129)	%dbSNP (b129)	Known Ts/Tv	Novel Ts/Tv	Overall Ts/Tv	%HM3 Rediscovery
PASS	17,490,556	6,768,218	38.7	2.12	2.07	2.09	90.164
FAIL	3,138,231	316,178	10.0	1.34	0.81	0.85	0.410
TOTAL	20,628,787	7,084,396	34.3	2.07	1.68	1.80	90.574

- Notice that the “failed” SNPs have:
  - lower fraction of previously known SNPs
  - Lower transition transversion ratio

# What is the right amount of filtering?

- For a population genetics study?
- For a disease association study?
  
- What might be pitfalls of too much filtering?
- What might be pitfalls of too little filtering?

# Rare variant test for high LDL vs low LDL

