

# **ESP6900 Call Set Overview**

March 28<sup>th</sup>, 2012

Goo Jun, Hyun Min Kang and Gonçalo Abecasis

Mark Rieder and Mark dePristo & teams

Leslie Lange and Paul Auer

# Overview of ESP6900 data set

- 6,916 samples included in calling
  - 5,407 QC-PASS samples from ESP5400
  - 1,509 new sampled added for this analysis
  - 70 TB of data coming in, 300GB VCF coming out
- ~100 TB of aligned sequence data
  - Avg. % consensus target bases covered per sample
    - $\geq 1x$  : 98%
    - $\geq 5x$  : 96%
    - $\geq 10x$  : 93%
    - $\geq 20x$  : 87%
    - $\geq 30x$  : 81%
- ~8 weeks of computational time for variant calling

# ESP6900 Sample Distribution

Category	Sample Distribution
By Target	BI (3,110), UW'08 (723), UW'09 (755), UW'10 (2,328)
By Cohort	ARIC(844), CARDIA(206), CHS (215), COPDGene(259), FHS(467), JHS(413), LHS(337), MESA(409), SARP(191), WHI(1,922), <i>UNKNOWN-TO-US (1,653)</i>
By Sex	Male (2,667), Female (4,249)
By Population	European (4,552) & African (2,364) Americans

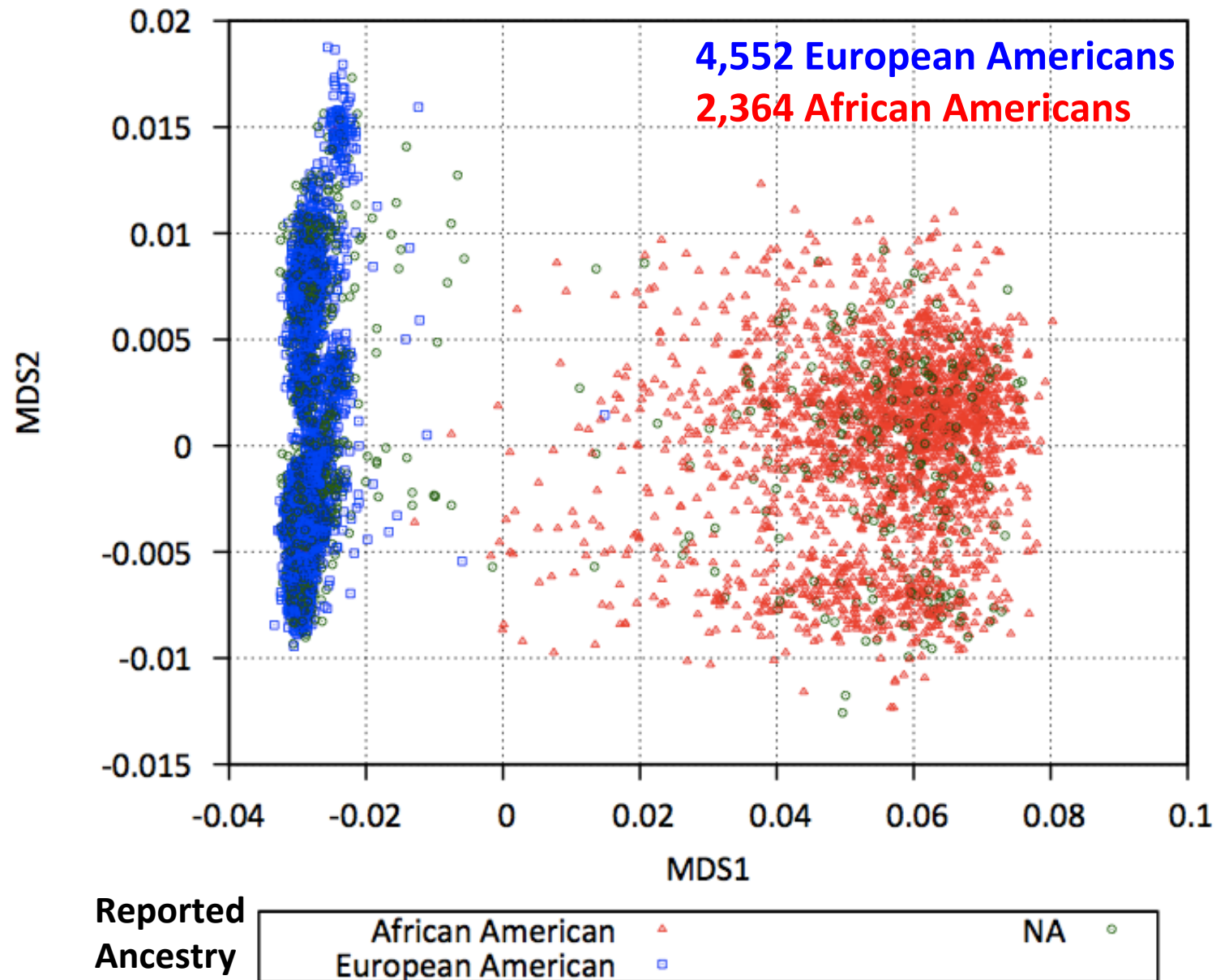
---

# Overview of ESP6900 SNP call sets

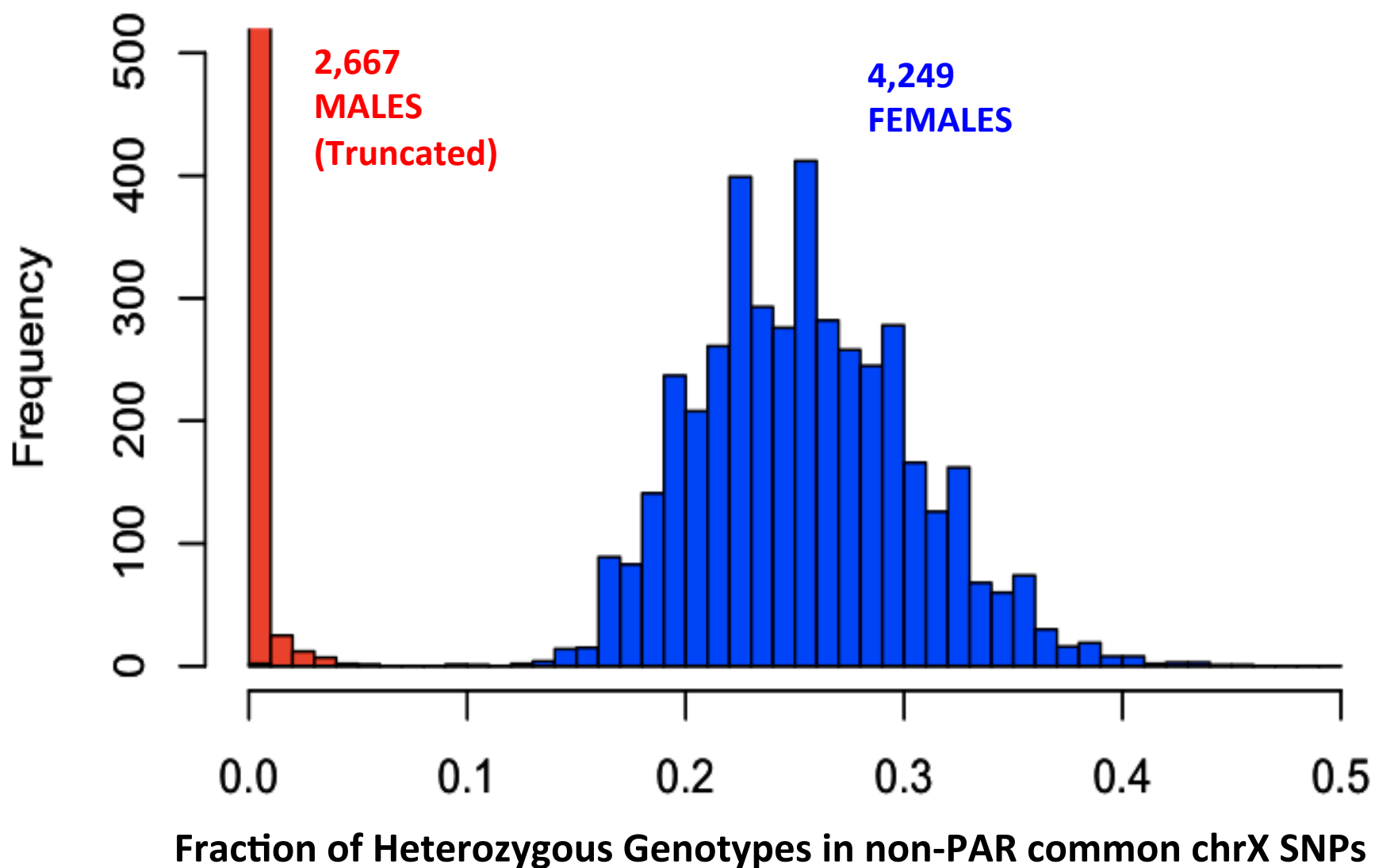
CATEGORY	#SNPs	% in dbSNP 129	% in 1000G Phase1	Known Ts/Tv	Novel Ts/Tv	%HM3 (ASW) Detected
UNFILTERED	2.10M	8.5	25.6	2.82	2.62	97.2
<b>PASS</b>	<b>1.92M</b>	<b>8.6</b>	<b>26.7</b>	<b>2.94</b>	<b>2.83</b>	<b>97.0</b>
<b>SILENT</b>	<b>433k</b>	<b>11.3</b>	<b>29.2</b>	<b>5.58</b>	<b>5.66</b>	<b>99.6</b>
<b>MISSENSE</b>	<b>729k</b>	<b>6.7</b>	<b>21.4</b>	<b>2.36</b>	<b>2.43</b>	<b>99.3</b>
<b>NONSENSE</b>	<b>19.6k</b>	<b>3.0</b>	<b>13.2</b>	<b>2.13</b>	<b>2.23</b>	<b>97.6</b>

- Variants are called within target $\pm$ 50bp
- Filters applied : SVM + nearby indels + strand bias filter for rare variants
- %HM3 sensitivity were evaluated in consensus target  $\pm$  50bp
- INDEL calling is in-progress : See Mark DePristo's presentation

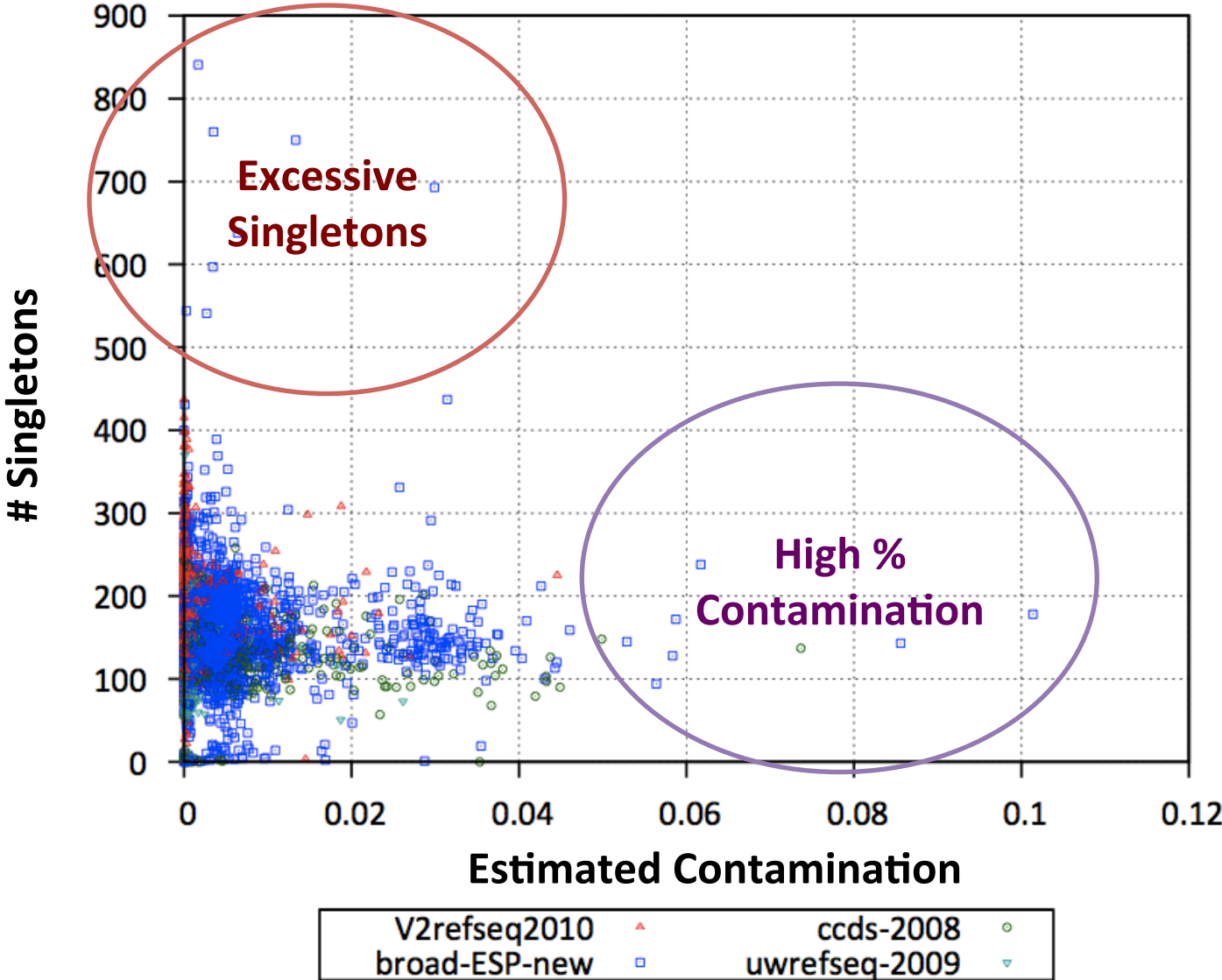
# MDS plot shows expected population clusters



# Sex is verified by chrX genotype likelihoods



# Outliers detected by various QC metrics



# Summary of sample level QC

# BEFORE QC	QC Applied	# AFTER QC
6,916	Sex Mismatch (7)	6,909
6,909	Excessive # of singletons (12)	6,897
6,896	>5% contamination (7)	6,890
6,890	>10% missing genotypes (7)	<b>6,883</b>
<b># samples included in the VCF release = 6,883</b>		
6,883	Duplicated Samples (48)	6,883
6,835	Related Samples (95)	<b>6,740</b>
<b>Preliminary # samples available for trait analysis = 6,740</b>		



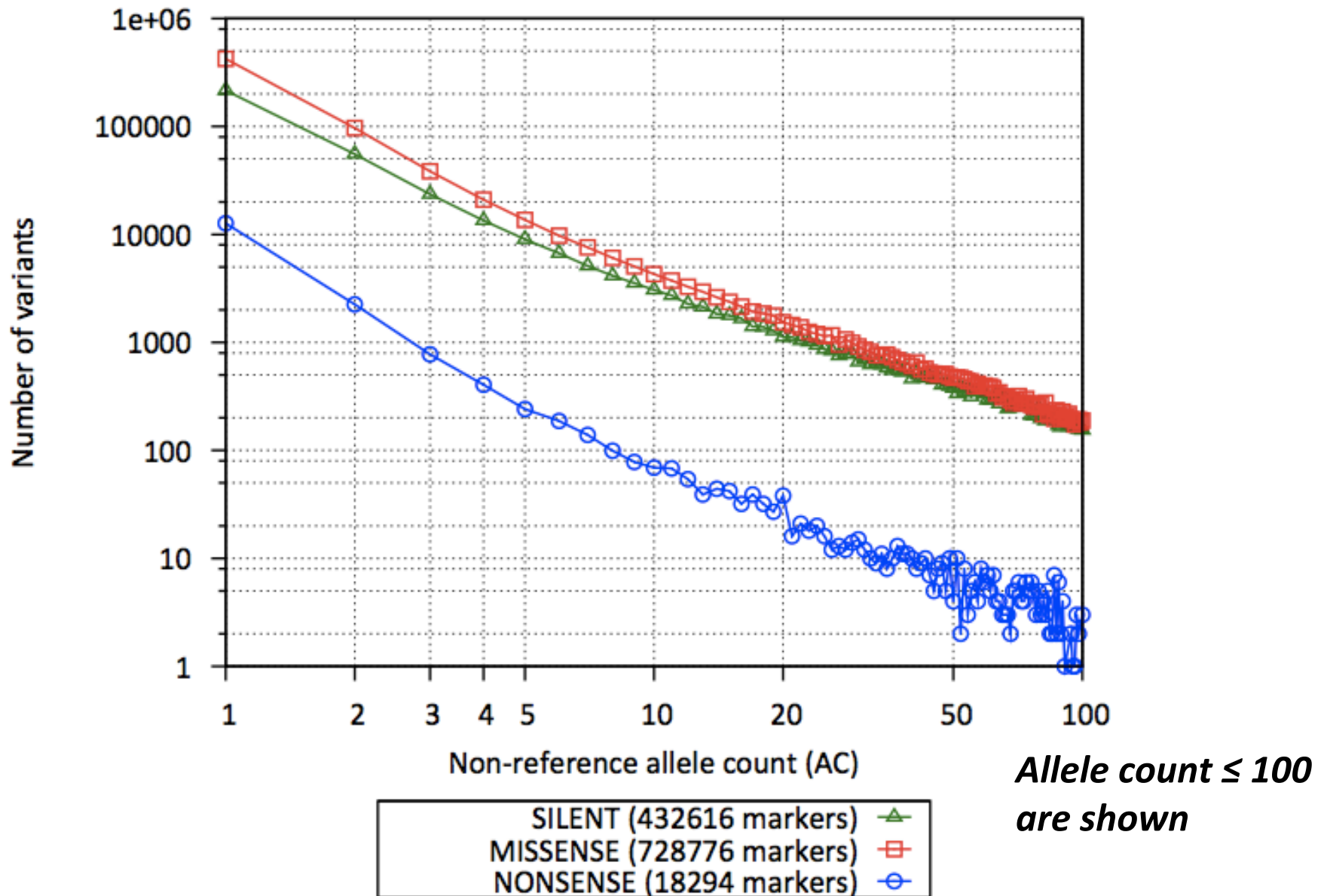
# Most Variants Are Rare

CATEGORY	#SNPs	SINGLE- TONs	DOUBLE- TONs	AF<0.1%	AF<0.5%	AF<5%
PASS	1.92M	53%	13%	82%	89%	96%
SILENT	433k	50%	13%	80%	88%	95%
MISSENSE	729k	58%	13%	87%	93%	97%
NONSENSE	19.6k	69%	12%	93%	97%	99%

# Singletons per individual :

**29 SILENT, 59 MISSENSE, 2 NONSENSE**

# Allele Frequency Spectrum

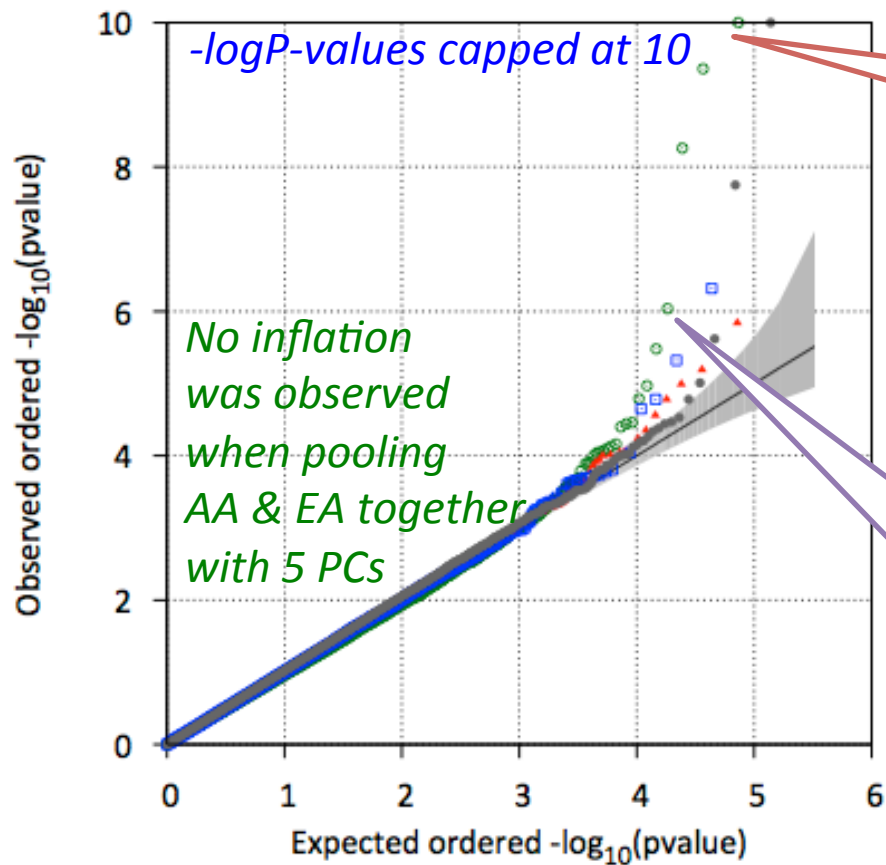


# How many variants per gene?

CATEGORY	SILENT	MISSENSE	NONSENSE
<b>AVG # SNPs per GENE</b>	<b>23.5</b>	<b>39.3</b>	<b>2.0</b>
<b>% GENES with <math>\geq 1</math> SNPs</b>	<b>98.9%</b>	<b>99.6%</b>	<b>48.7%</b>
<b>% GENES with <math>\geq 5</math> SNPs</b>	<b>91.5%</b>	<b>95.3%</b>	<b>3.1%</b>
<b>% GENES with <math>\geq 10</math> SNPs</b>	<b>74.7%</b>	<b>87.0%</b>	<b>0.1%</b>
<b>% GENES with <math>\geq 20</math> SNPs</b>	<b>42.6%</b>	<b>67.0%</b>	<b>0.0%</b>
<b>% GENES with <math>\geq 50</math> SNPs</b>	<b>9.0%</b>	<b>24.5%</b>	<b>0.0%</b>
<b>% GENES with <math>\geq 100</math> SNPs</b>	<b>1.7%</b>	<b>5.7%</b>	<b>0.0%</b>

- Gene length and mutation rate plays major role in statistical power given effect size
- Results are based on GenCodeV7 annotations

# Preliminary single variant analysis on LDL



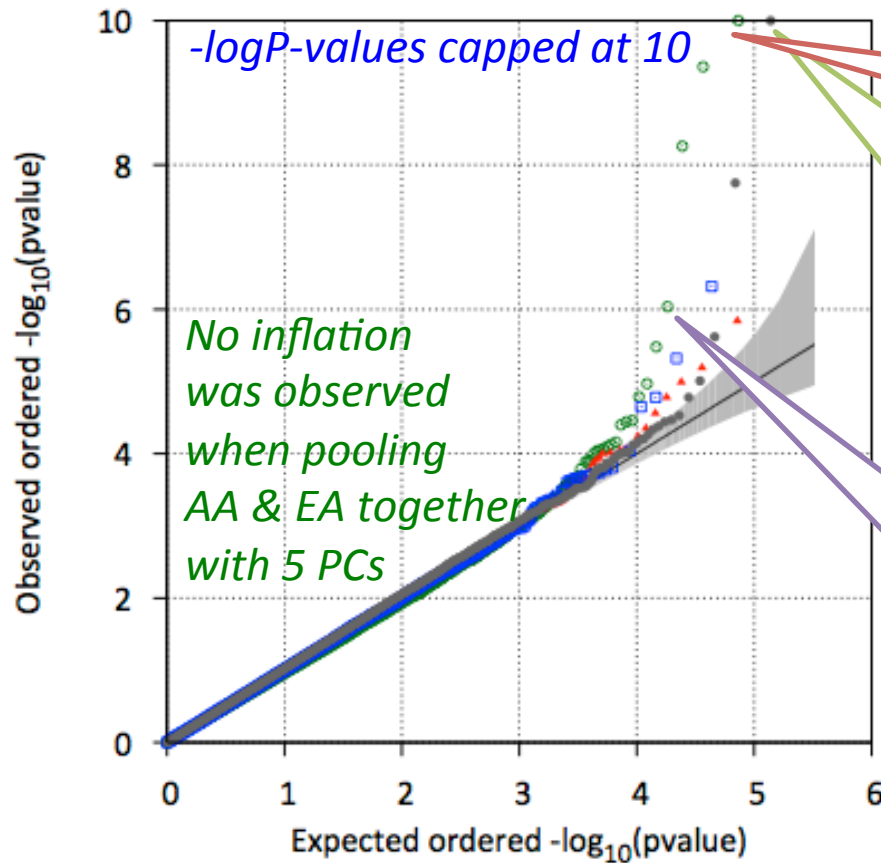
**TOMM40 (near APOE)**  
Synonymous SNP  
MAF 0.065  
p-value :  $4.2 \times 10^{-24}$

**APOB**  
Missense SNP  
MAF 0.21  
p-value :  $9.0 \times 10^{-7}$

MAF [.05,.5] (73306)	○
MAF [.01,.05] (71916)	▲
MAF [.005,.01] (43622)	□
MAF [.001,.005] (138249)	●

*N=3,342 individuals with phenotypes*

# Preliminary single variant analysis on LDL



MAF [.05,.5] (73306)	○
MAF [.01,.05] (71916)	▲
MAF [.005,.01] (43622)	□
MAF [.001,.005] (138249)	●

**TOMM40 (near APOE)**  
 Synonymous SNP  
 MAF 0.065  
 p-value :  $4.2 \times 10^{-24}$

**PCSK9**  
 Nonsense SNP 1      Nonsense SNP 2  
 MAF 0.005              MAF 0.002  
 p-value :  $7.2 \times 10^{-11}$       p-value :  $1.8 \times 10^{-8}$

**APOB**  
 Missense SNP  
 MAF 0.21  
 p-value :  $9.0 \times 10^{-7}$

*N=3,342 individuals with phenotypes*

# Summary

- ESP6900 SNP call sets are produced
  - Harmonized across different centers
  - Vast majority of them are novel rare variants
  - ~730k missense, ~430k silent, ~20k nonsense SNPs
- Rare variant association with larger sample size
  - Signals for known SNPs increases by orders of magnitude
  - Nonsense SNPs at  $MAF < .5\%$  are genome-wide significant
  - Integrated and harmonized phenotype analyses likely improve statistical power

# Acknowledgements

- @Michigan
  - Goo Jun
  - Hyun Min Kang
  - Gonçalo Abecasis
  - Youna Hu
  - Cristen Willer
  - Chenyi Xue
- @Broad
  - Mark DePristo
  - Stacey Gabriel
- @UW
  - Mark Rieder
  - Debbie Nickerson
- Leslie Lange @ UNC
- Paul Auer @ FHCRC
- Suzanne Leal @ BCM