

# The 1000 Genomes Project Tutorial

ICHG 2011

Montreal, Quebec, Canada

October 13, 2011





# 1000 Genomes

A Deep Catalog of Human Genetic Variation

- International project to construct a foundational data set for human genetics
  - Discover virtually all common human variations by investigating many genomes at the base pair level
  - Consortium with multiple centers, platforms, funders
- Aims
  - Discover population level human genetic variations of all types (95% of variation > 1% frequency)
  - Define haplotype structure in the human genome
  - Develop sequence analysis methods, tools, and other reagents that can be transferred to other sequencing projects

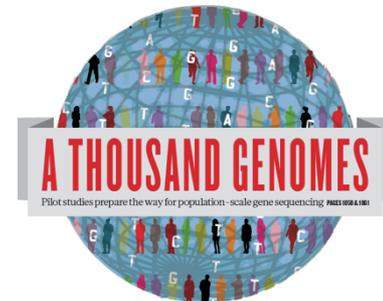
# Agenda

<b>Time</b>	<b>Topic</b>	<b>Presenter</b>	<b>Presenter affiliation</b>
7:30	Description of 1000 Genomes data	Gabor Marth, D.Sc.	Boston College, Boston, MA
7:55	How to access the data	Paul Flicek, D.Sc.	EMBL European Bioinformatics Inst., Hinxton, Cambridge, UK
8:20	Lessons in variant calling and genotyping	Hyun Min Kang, Ph.D.	Univ. of Michigan, Ann Arbor, MI
8:40	Structural variants	Ryan Mills, Ph.D.	Brigham and Women's Hospital, Boston, MA
9:00	Imputation in GWAS studies	Bryan Howie, Ph.D.	Univ. of Chicago, Chicago, IL
9:20	Q&A	-	-

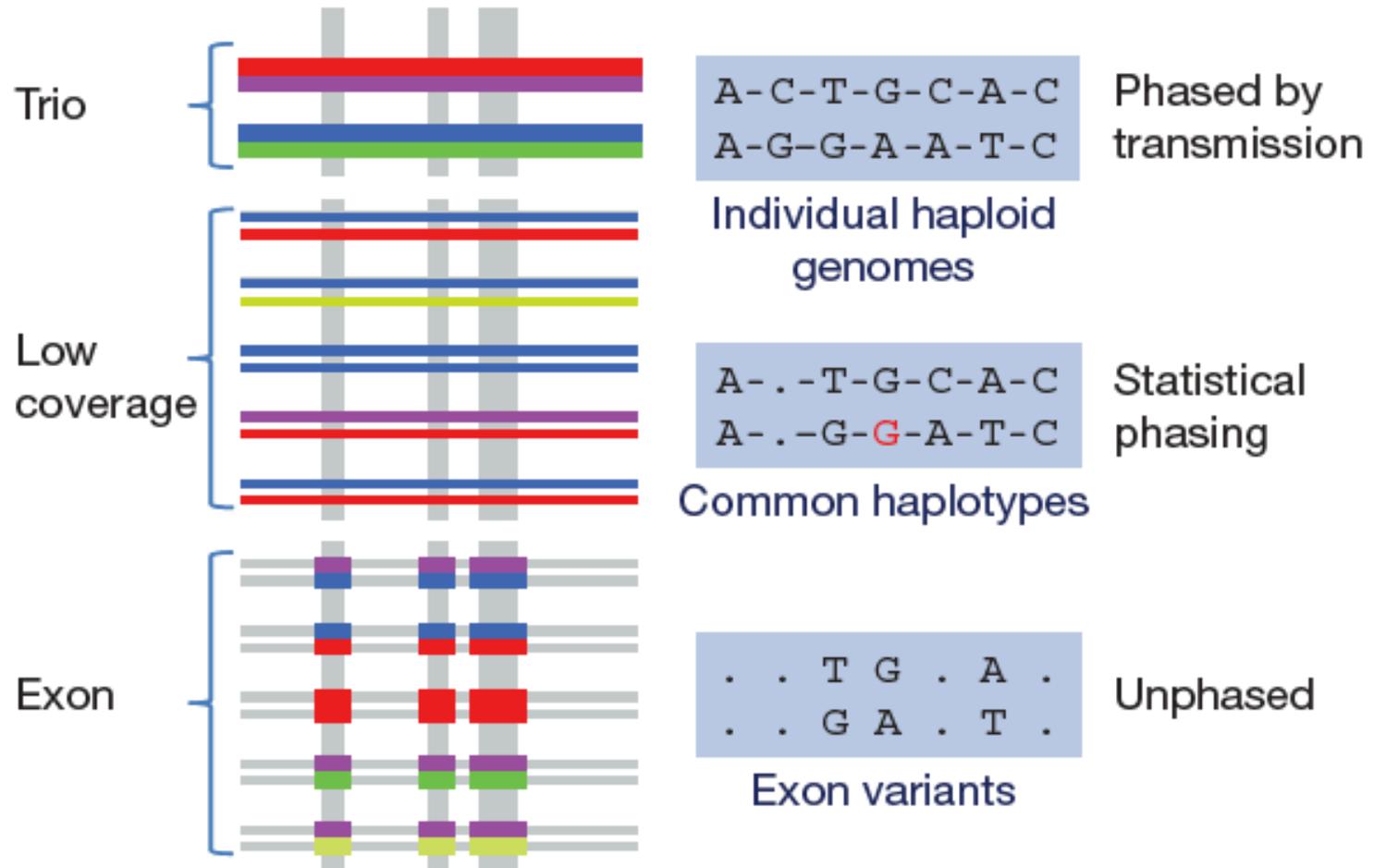
# The 1000 Genomes Project Datasets

**Gabor T. Marth**  
**Boston College Biology**  
**Department**

1000 Genomes Project Tutorial  
Montreal, Quebec, Canada  
October 13, 2011



# 3 pilot coverage strategies



# Pilot results published

## ARTICLE

doi:10.1038/nature09534

# A map of human genome variation from population-scale sequencing

The 1000 Genomes Project Consortium\*

Marth *et al.* *Genome Biology* 2011, 12:R84  
http://genomebiology.com/2011/12/9/R84



OPEN ACCESS Freely available online

PLoS GENETICS

RESEARCH

Open Access

## A Comprehensive Map of Mobile Element Insertion Polymorphisms in Humans

Chip Stewart<sup>1</sup>\*, Deniz Kural<sup>1</sup>\*, Michael P. Strömberg<sup>1</sup>\*, Jerilyn A. Walker<sup>2</sup>, Miriam K. Konkel<sup>2</sup>, Adrian M. Stütz<sup>3</sup>, Alexander E. Urban<sup>4</sup>, Fabian Grubert<sup>4</sup>, Hugo Y. K. Lam<sup>4</sup>, Wan-Ping Lee<sup>1</sup>, Michele Busby<sup>1</sup>, Amit R. Indap<sup>1</sup>, Erik Garrison<sup>1</sup>, Chad Huff<sup>2</sup>, Jinchuan Xing<sup>5</sup>, Michael P. Snyder<sup>4</sup>, Lynn B. Jorde<sup>4</sup>, Mark A. Batzer<sup>2</sup>, Jan O. Korbel<sup>3</sup>, Gabor T. Marth<sup>1</sup>\*, 1000 Genomes Project<sup>1</sup>

<sup>1</sup>Department of Biology, Boston College, Chestnut Hill, Massachusetts, United States of America, <sup>2</sup>Department of Biological Sciences, Louisiana State University, Baton Rouge, Louisiana, United States of America, <sup>3</sup>Genome Biology Unit, European Molecular Biology Laboratory, Heidelberg, Germany, <sup>4</sup>Department of Genetics, Stanford University, Stanford, California, United States of America, <sup>5</sup>Department of Human Genetics, Eccles Institute of Human Genetics, University of Utah, Salt Lake City, Utah, United States of America

## The functional spectrum of low-frequency coding variation

Gabor T Marth<sup>1</sup>\*, Fuli Yu<sup>2</sup>†, Amit R Indap<sup>1</sup>†, Kiran Garimella<sup>3</sup>†, Simon Gravel<sup>4</sup>†, Wen Fung Leong<sup>1</sup>†, Chris Tyler-Smith<sup>5</sup>†, Matthew Bainbridge<sup>2</sup>, Tom Blackwell<sup>6</sup>, Xiangqun Zheng-Bradley<sup>7</sup>, Yuan Chen<sup>5</sup>, Danny Challis<sup>2</sup>, Laura Clarke<sup>7</sup>, Edward V Ball<sup>8</sup>, Kristian Cibulskis<sup>3</sup>, David N Cooper<sup>8</sup>, Bob Fulton<sup>9</sup>, Chris Hart<sup>3</sup>, Dan Koboldt<sup>9</sup>, Donna Muzny<sup>4</sup>, Richard Smith<sup>7</sup>, Carrie Sougnez<sup>3</sup>, Chip Stewart<sup>1</sup>, Alistair Ward<sup>1</sup>, Jin Yu<sup>2</sup>, Yali Xue<sup>5</sup>, David Altshuler<sup>3</sup>, Carlos D Bustamante<sup>4</sup>, Andrew G Clark<sup>10</sup>, Mark Daly<sup>3</sup>, Mark DePristo<sup>3</sup>, Paul Flicek<sup>7</sup>, Stacey Gabriel<sup>3</sup>, Elaine Mardis<sup>9</sup>, Aarno Palotie<sup>5</sup>, Richard Gibbs<sup>2</sup> and the 1000 Genomes Project

## ARTICLE

doi:10.1038/nature09708

## Mapping copy number variation by population-scale genome sequencing

Ryan E. Mills<sup>1</sup>\*, Klaudia Walter<sup>2</sup>\*, Chip Stewart<sup>3</sup>\*, Robert E. Handsaker<sup>4</sup>\*, Ken Chen<sup>5</sup>\*, Can Alkan<sup>6,7\*</sup>, Alexej Abyzov<sup>8\*</sup>, Seungtae Chris Yoon<sup>9\*</sup>, Kai Ye<sup>10\*</sup>, R. Keira Cheetham<sup>1</sup>, Asif Chinwalla<sup>5</sup>, Donald F. Conrad<sup>2</sup>, Yutao Fu<sup>12</sup>, Fabian Grubert<sup>13</sup>, Iman Hajirasouliha<sup>14</sup>, Fereydoon Hormozdliari<sup>14</sup>, Lilla M. Iakoucheva<sup>15</sup>, Zamin Iqbal<sup>16</sup>, Shuli Kang<sup>16</sup>, Jeffrey M. Kidd<sup>16</sup>, Miriam K. Konkel<sup>17</sup>, Joshua Korn<sup>18</sup>, Ekta Khurana<sup>8,18</sup>, Deniz Kural<sup>1</sup>, Hugo Y. K. Lam<sup>13</sup>, Jing Leng<sup>8</sup>, Ruiqiang Li<sup>19</sup>, Yingrui Li<sup>19</sup>, Chang-Yun Lin<sup>20</sup>, Ruibang Luo<sup>19</sup>, Xinmeng Jasmine Mu<sup>8</sup>, James Nemesh<sup>1</sup>, Heather E. Peckham<sup>12</sup>, Tobias Rausch<sup>21</sup>, Aylwyn Scally<sup>2</sup>, Xinghua Shi<sup>1</sup>, Michael P. Stromberg<sup>3</sup>, Adrian M. Stütz<sup>3</sup>, Alexander Ekekehart Urban<sup>13,27</sup>, Jerilyn A. Walker<sup>17</sup>, Jiantao Wu<sup>1</sup>, Yujun Zhang<sup>2</sup>, Zhengdong D. Zhang<sup>8</sup>, Mark A. Batzer<sup>17</sup>, Li Ding<sup>22</sup>, Gabor T. Marth<sup>1</sup>, Gil McVean<sup>23</sup>, Jonathan Sebat<sup>3</sup>, Michael Snyder<sup>13</sup>, Jun Wang<sup>19,24</sup>, Kenny Ye<sup>20</sup>, Evan E. Eichler<sup>25</sup>, Mark B. Gerstein<sup>8,18,25</sup>, Matthew E. Hurles<sup>2</sup>, Charles Lee<sup>1</sup>, Steven A. McCarroll<sup>1,26</sup>, Jan O. Korbel<sup>21</sup> & 1000 Genomes Project

## Demographic history and rare allele sharing among human populations

Simon Gravel<sup>a</sup>, Brenna M. Henn<sup>a</sup>, Ryan N. Gutenkunst<sup>b</sup>, Amit R. Indap<sup>c</sup>, Gabor T. Marth<sup>c</sup>, Andrew G. Clark<sup>d</sup>, Fuli Yu<sup>e</sup>, Richard A. Gibbs<sup>e</sup>, The 1000 Genomes Project<sup>e</sup>, and Carlos D. Bustamante<sup>e,1</sup>

<sup>a</sup>Department of Genetics, Stanford University School of Medicine, Stanford, CA 94305-5120; <sup>b</sup>Department of Molecular and Cellular Biology, University of Arizona, Tucson, AZ 85721; <sup>c</sup>Department of Biology, Boston College, Chestnut Hill, MA 02467; <sup>d</sup>Department of Molecular Biology and Genetics, Cornell University, Ithaca, NY 14853; and <sup>e</sup>Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX 77030

Edited by Michael Lynch, Indiana University, Bloomington, IN, and approved June 3, 2011 (received for review December 24, 2010)

# Finalized project design

- Based on the result of the pilot project, we decided to collect data on 2,500 samples from 5 continental groupings
  - Whole-genome low coverage data (>4x)
  - Full exome data at deep coverage (>50x)
  - A number of deep coverage genomes to be sequenced, with details to be decided
  - Hi-density genotyping at subsets of sites
- Moved from the Pilot into Phase 1 of the project

# Phase I (1,150)

# Phase II (1,721)

# Phase III (2,500)

CDX 17S



CLM (70T); DNA from LCL



CHS (100T); DNA from LCL



PUR (70T); DNA from Blood



FIN (100S); DNA from LCL



GBR (96/100S); DNA from LCL



IBS (84/100T); DNA from LCL



GWD



GWD



GWD



GWD (target - 100T); DNA from LCL



CDX (100S); DNA: 17 DNA from Bld, 83 from LCL



KHV (82/100) - 15 trios; DNA Bld



45      99 (29T)      23 (7T)

ACB (28/79T) - 14 trios; DNA Bld



13 26 20 9 26 39 27 26 22

PEL (70T); DNA from Blood



3



1



16 (8T)



PJL (target - 100T); DNA from Blood



15      6      6      195

GWD



GWD



GWD



GWD (target - 100T); DNA from LCL



12      15      15      270

GIH vs. Sindhi (target - 100T)



Tamil (target - 100T)



Sri Lankan (target - 100T)



Bengalee (target - 100T)



Nigeria (target - 100T); DNA from LCL



Sierra Leone (target - 100T); DNA from LCL



MAB (target - 100T); DNA from LCL



AJM (target - 80T); DNA from Bld



April 2009    June 2009    Aug 2009    Oct 2009    Dec 2009    Feb 2010    April 2010    June 2010    Aug 2010    Oct 2010    Dec 2010    Feb 2011    April 2011    June 2011    Aug 2011    Oct 2011    Dec 2011    Feb 2012    April 2012

# Phase I data

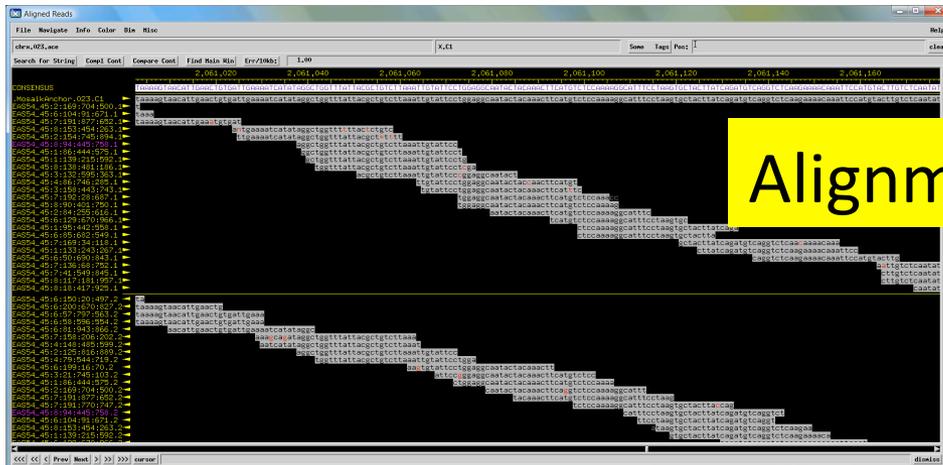
- Samples from 14 populations: ASW, CEU, CHB, CHS, CLM, FIN, GBR, IBS, JPT, LWK, MXL, PRU, TSI, YRI

Dataset	Low coverage whole genome	Deep coverage whole exome
# samples	1,094	1,128
Sequencing technologies	Illumina, SOLiD, 454	Illumina, SOLiD
Primary alignments (BAMs)	BWA, BFAST	MOSAIK, BFAST
Second alignments (BAMs)	MOSAIK	BWA, MOSAIK
Read coverage	4-8X per sample	≥70% of targets with ≥20X coverage in every sample

# Raw data & read alignment delivery

```
@IL11_266:1:1:395:231/1  
CCAACCACAACACAAAAACACAAGCAACACACAC  
+  
@AAAAA?<>@@>?:475;A6?384,>5  
@IL11_266:1:1:399:301/1  
CAAAAAAAGAAAGTACGAGATACGACACATCAC  
+  
;@AAAA>5;>@C67'&2?&7<&7&@1/1408=19::
```

Reads: FASTQ

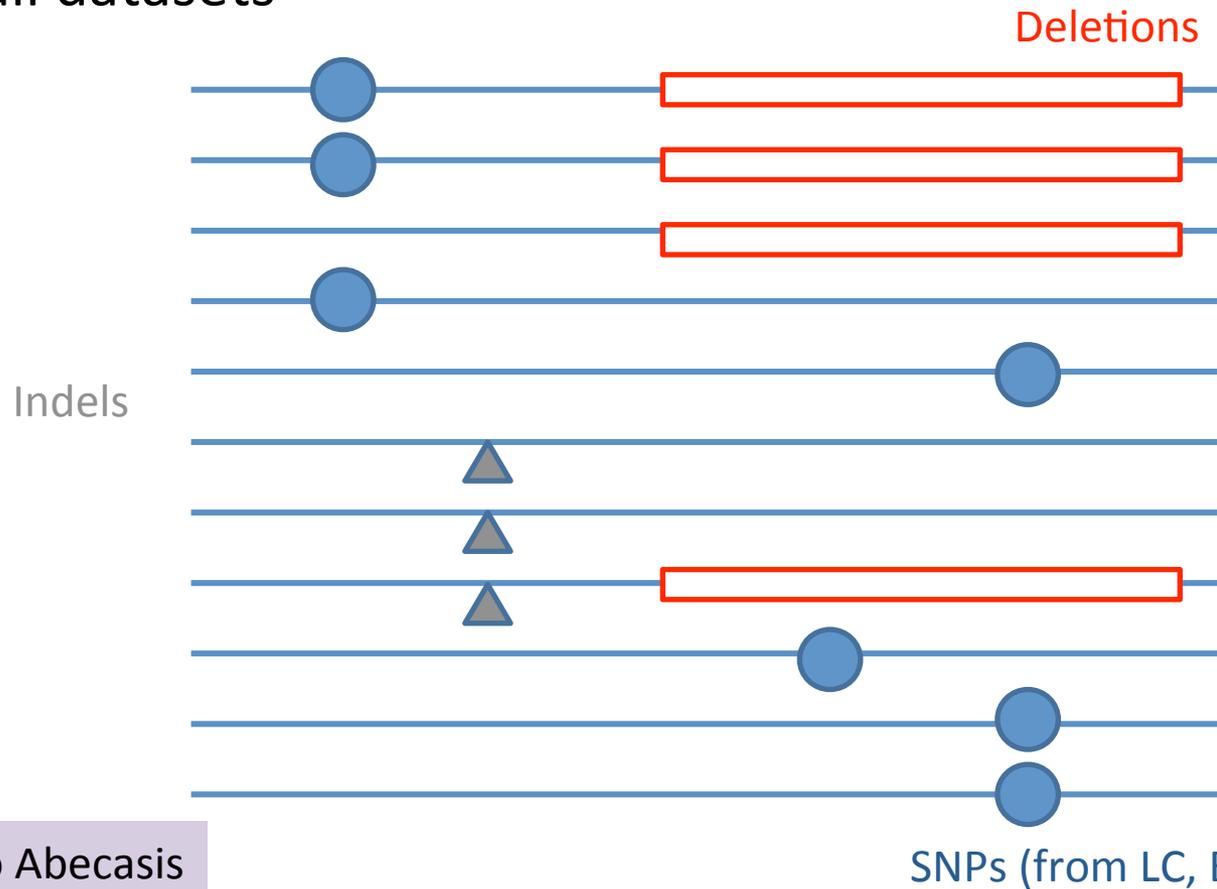


Alignments: BAM

<ftp://ftp.1000genomes.ebi.ac.uk>

# Phase 1 analysis goal: an **integrated view of human variations**

- Reconstruct haplotypes including all variant types, using all datasets



# Pipelines for data processing and variant calling

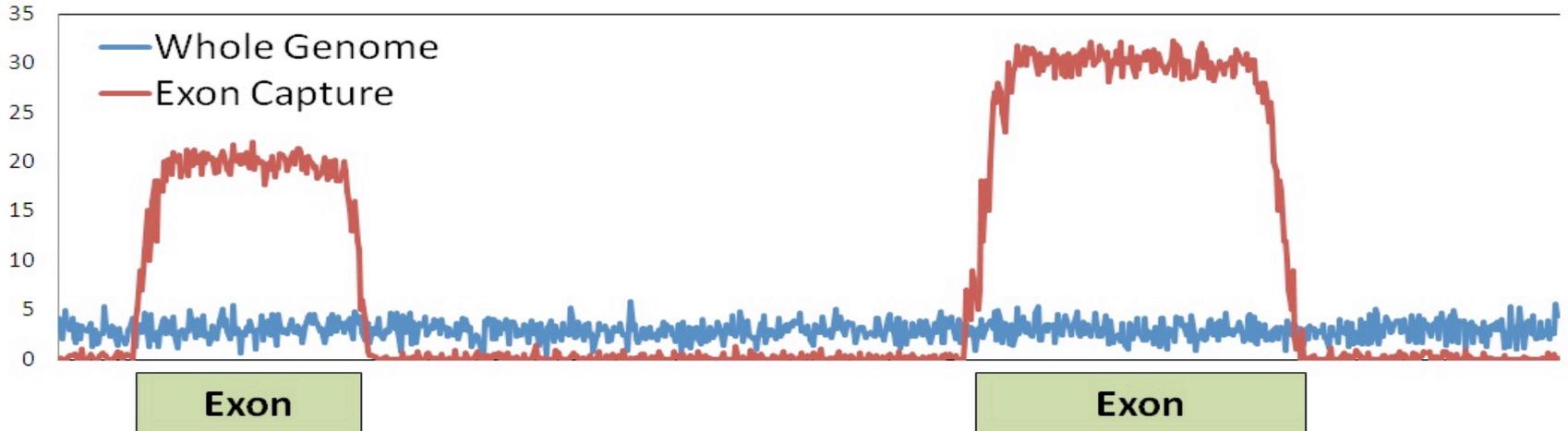
- Tens of analysis groups have contributed
- Individual pipelines and component tools vary
- Typical main steps:
  - Read mapping
  - Duplicate filtering
  - Base quality value recalibration
  - INDEL realignment
  - Variant calling (sites)
  - Sample genotype calling (sometime part of variant calling)
  - Variant filtering / call set refinement
  - Variant reporting

# SNPs

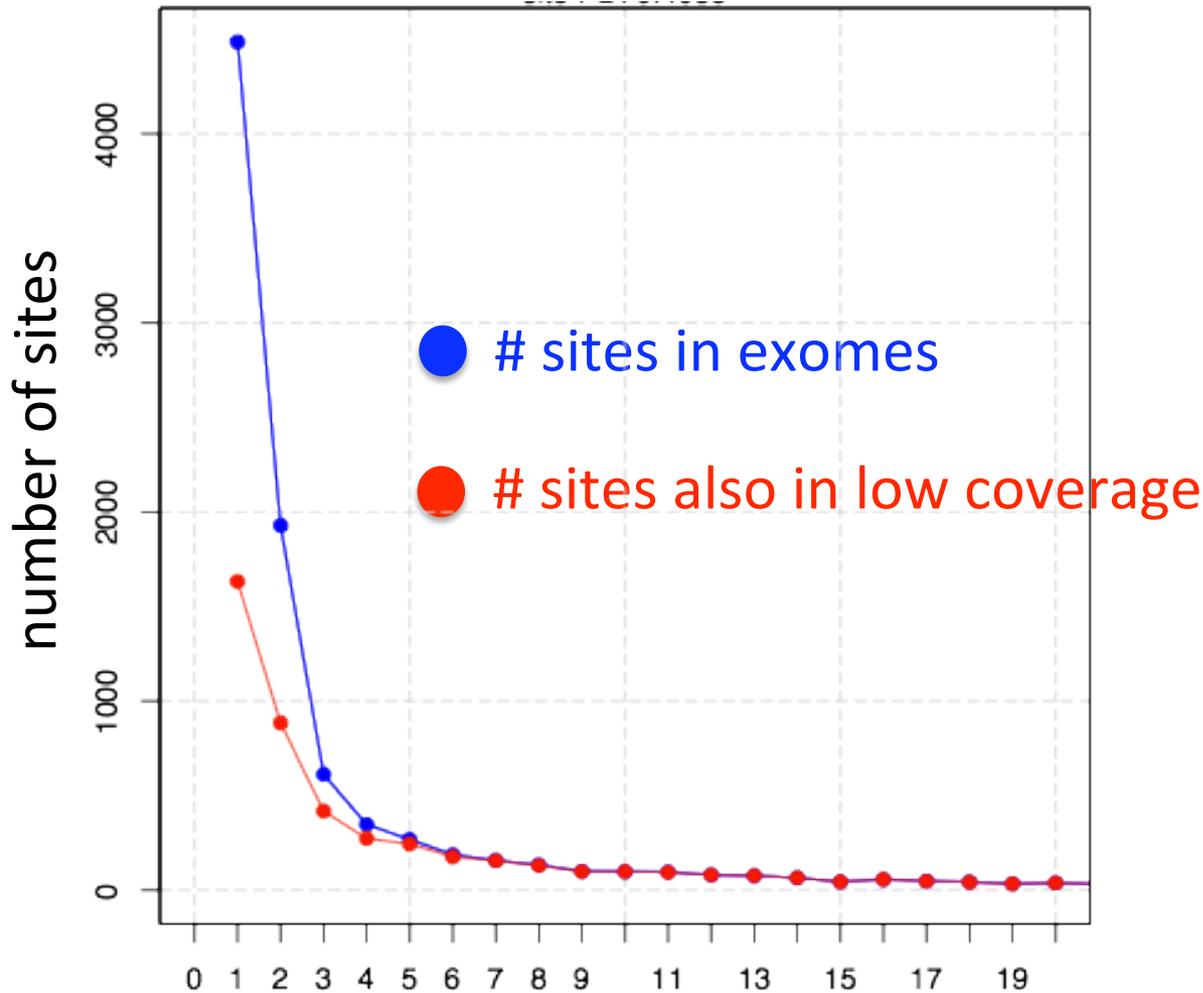


# SNP calls

Dataset	Contributing datasets	Consensus method	#SNPs	# Novel SNPs	Novel Ts/Tv	%ONMI poly (sensitivity)	%OMNI mono (FDR)
Low coverage	BC, BCM, BI, NCBI, UM	VQSR	37.9M	29.65M	2.16	98.4	1.80
Exome/Illumina	BC, BCM, BI, Cornell, UM	SVM	598K	468K	2.74	98.01	1.97
Exome/SOLiD	BC, BCM, UM	SVM	356K	243K	2.91	90.67	1.29



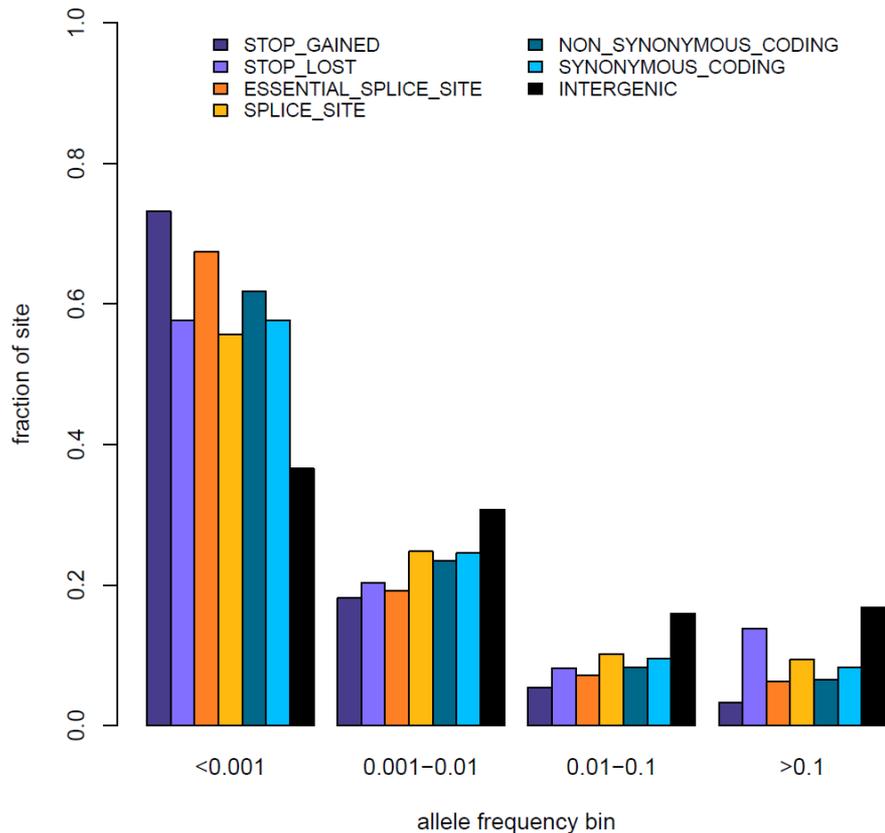
# Deep coverage exome data is more sensitive to low-frequency variants



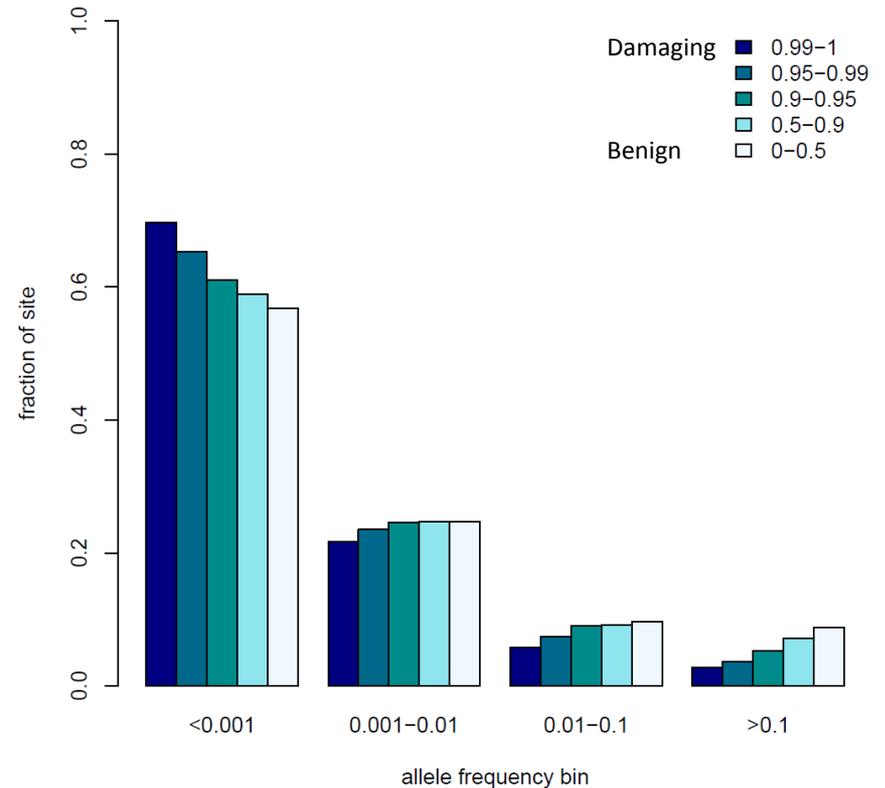
Allele count in 766 exomes (chr. 20, exons only)

# Newly discovered SNPs are mostly at low frequency and enriched for functional variants

## Functional category



## Non-synonymous: Condel score



# INDELS

```
tttatttagggctgagcaataatag
tttatttagggctgagcaataatag
tttatttagggctgagcaataatag
tttatttagggctgagc**taatagacg
      ttagggctgagcaataatagacg
            agggctgagc**taatagacg
                  agggctgagc**taatagacg
                        gctgagc**taatagacg
                              tgagc**taatagacg
                                    tgagc**taatagacg
                                          tgagcaataatagacg
                                                gagc**taatagacg
                                                      gagc**taatagacg
                                                            gagc**taatagacg
                                                                  agcaataatagacg
                                                                        gc**taatagacg
                                                                              taatagacg
                                                                                      agacg
                                                                                          gacg
```

```
GATTAGAATCGCAATTAAC
GATTAGAAT*GCAATTAAC
```

```
AGTTTCTCT***TTCTTACAG
AGTTTCTCTGCTTTCTTACAG
```

```
CGAATTAGA*****GCAA
CGAATTAGACTTAGAGCAA
```

```
TCTCAAAAAAAAAAAAAAAAAAAGTGT
TCTCAAAAAAAAAAAAAAAAAA*GTGT
```

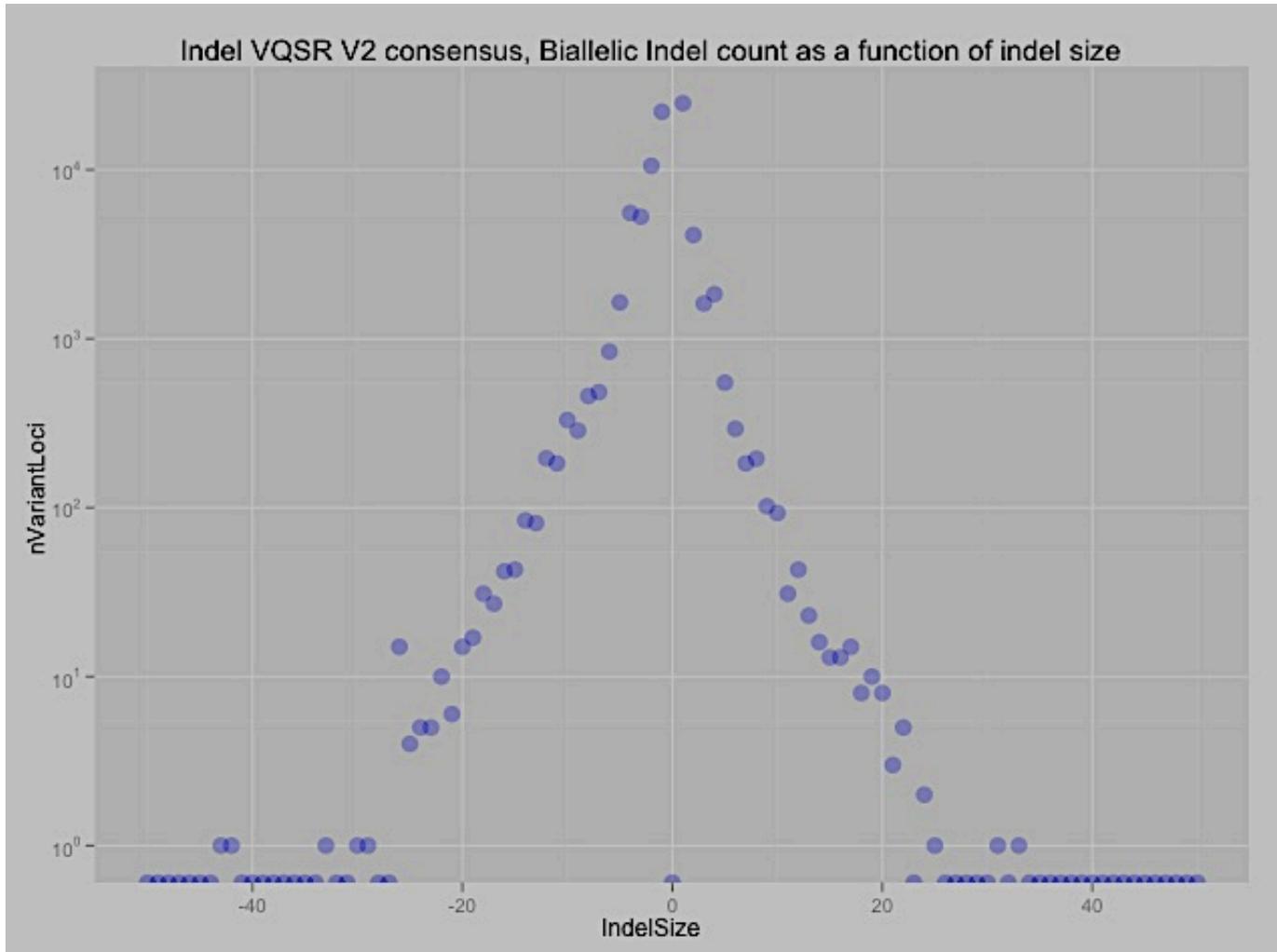
```
YAAA*****GA
YAAAAAAAAAAAAAAAAAAAAAGA
```

```
TGTGTGTGTGTGTGTGTATTTAAAAACTAGG
TGTGTGTGTGTGTGTG**TATTTAAAAACTAGG
```

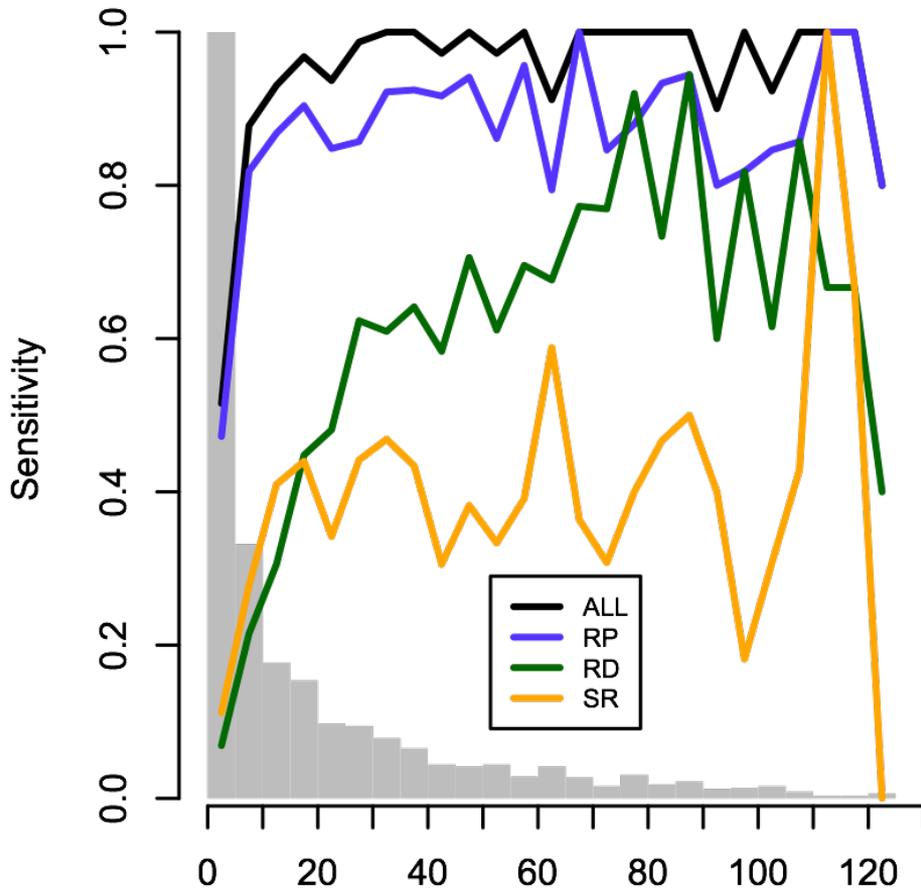
# INDEL calls

Dataset	Contributing datasets	Consensus method	#INDELs
Low coverage	BC, BI, DI, OX, SI	VQSR	5.5M
Exome/Illumina	BC, BCM, BI	N.A.	6.5 – 10.2K
Exome/SOLiD	BCM	N.A.	4.2 – 5.0K

# INDEL length



# Finding structural variants



- Discovery with a number of different methods
- Several types (e.g. deletions, tandem duplications, mobile element insertions) now detectable with high accuracy
- We are pulling in new types for the Phase I data (inversions, *de novo* insertions, translocations)

# SNP validations (low coverage data)

	Total	Polymorphic	Monomorphic	No Call	Confirmation Rate	Failure Rate
<b>All Sites</b>	<b>300</b>	<b>282</b>	<b>12</b>	<b>6</b>	<b>0.959</b>	<b>0.020</b>
<b>Called in Validation Samples</b>	<b>287</b>	<b>276</b>	<b>5</b>	<b>6</b>	<b>0.982</b>	<b>0.021</b>
Singletons	70	65	3	2	0.956	0.029
MAF<0.01*	134	131	2	1	0.985	0.007
0.01<MAF<0.05	33	33	0	0	1.000	0.000
MAF>0.05	50	47	0	3	1.000	0.060

\*Excludes singletons

Danny Challis, Eric Banks

# Genotypes are accurate

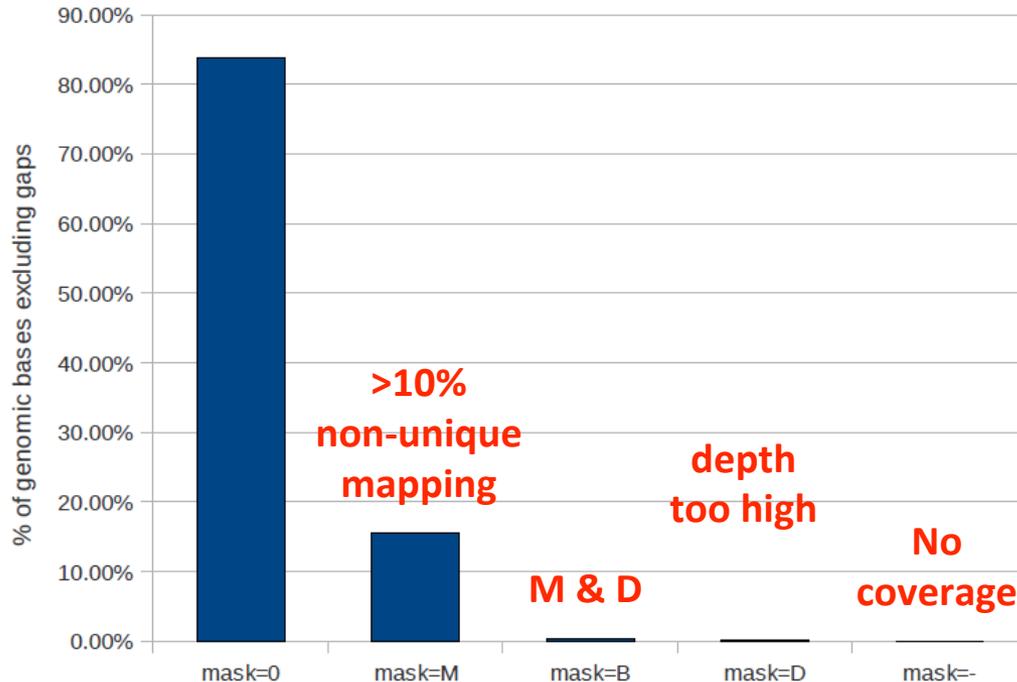
- Average low coverage depth is ~5x
- We obtain genotypes by sharing data between samples (using imputation-related methods)

<b>Genotype</b>	<b>HomRef</b>	<b>Het</b>	<b>HomAlt</b>	<b>Overall</b>
<b>Error rate</b>	0.16%	0.76%	0.39%	0.37%

- Genotypes are expected to be even more accurate after integration of multiple variant sources

# Accessible fraction of genome

Genomic coverage of mask types



- In the Pilot data, we found that >80% of the human genome reference was accessible for SNP variant calling
- We are currently re-evaluating this fraction for the Phase 1 data (which used longer reads)
- We are developing methods to estimate the fraction for other variants (especially INDELS)

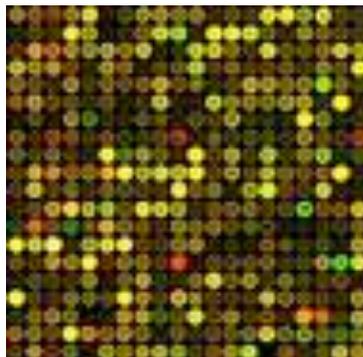
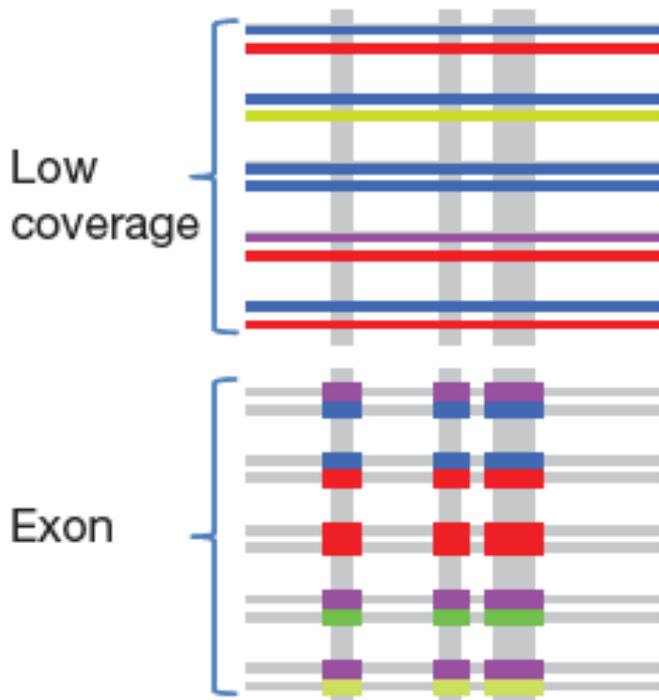
# Variant call delivery

Format: VCF

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	NA00001	NA00002
20	14370	rs6054257	G	A	29	0	NS=3;DP=14;AF=0.5;DB;H2	GT:GQ:DP:HQ	0 0:48:1:51,51	1 0:48:8:51,51
20	17330	.	T	A	3	q10	NS=3;DP=11;AF=0.017	GT:GQ:DP:HQ	0 0:49:3:58,50	0 1:3:5:65,3
20	1110696	rs6040355	A	G,T	67	0	NS=2;DP=10;AF=0.333,0.667;AA=T;DB	GT:GQ:DP:HQ	1 2:21:6:23,27	2 1:2:0:18,2
20	1230237	.	T	.	47	0	NS=3;DP=13;AA=T	GT:GQ:DP:HQ	0 0:54:7:56,60	0 0:48:4:51,51
20	1234567	microsat1	G	D4,IGA	50	0	NS=3;DP=9;AA=G	GT:GQ:DP	0/1:35:4	0/2:17:2

<ftp://ftp.1000genomes.ebi.ac.uk>

# Datasets & variant types

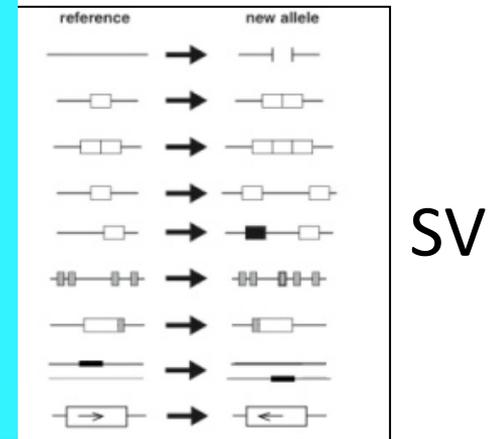


CTGAG  
ATGAG

SNP

CCCTGAG  
--TGAG

INDEL



# Data delivery

**1000 Genomes Pilot**  
A Deep Catalog of Human Genetic Variation

Tools | Help

**Search 1000 Genomes**

e.g. gene BRCA2 or Chromosome 6:133017695-133161157

**Start Browsing 1000 Genomes data**

-  [Browse Human](#) →  
NCBI 36
- [Transcript SNP view](#) →  
View the consequences of sequence variation at the level of each transcript in the genome.
- [Sequence Alignment View](#) →  
Shows read-depth data alongside SNPs

**Pilot Browser**

based on the full pilot project data described in [A map of human genome variation from population-scale sequencing](#), Nature 467, 1061.1073.

Please see [www.1000genomes.org](http://www.1000genomes.org) for more information about the data presented here and instructions for downloading the complete data set.

- [View sample data](#)

**The 1000 Genomes Browser**

Ensembl-based browser provides access to 1000genomes data

This browser represents the variant set analysed as part of [A map of human genome variation from population-scale sequencing](#), Nature 467, 1061.1073. The data behind this browser can be found on [the 1000 Genomes ftp site](#). This data can also be found in Ensembl and UCSC.

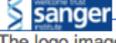
**Links**

-  [1000 Genomes](#) →  
More information about the 1000 Genomes Project on the 1000 genomes main site.

The 1000 Genomes Project is an international collaborative project described at [www.1000genomes.org](http://www.1000genomes.org).

The 1000 Genomes Browser is based on Ensembl web code.

[Ensembl](#) is a joint project of EMBL-EBI  and the [Wellcome Trust Sanger Institute](#)

 The logo image courtesy of [Andy Martin](#)

1000 Genomes Pilot release 7 - May 2011 © [EBI](#)

[About 1000 Genomes](#) | [Contact Us](#) | [Help](#)

Presentation on data access by Paul Flicek

# The 1000GP is a driver for method and tool development

- New data formats (SAM/BAM, VCF) developed by the 1000GP are now adopted by the entire genomics community
- Tools (read mappers e.g. BWA, MOSAIK, etc; variant callers including those for SVs)
- Data processing protocols (BQ recalibration, duplicate read removal, etc.)
- Imputation and haplotype phasing methods

# Tools for analyzing & manipulating 1000G data



Alignments: SAM/BAM

- samtools: <http://samtools.sourceforge.net/>
- BamTools: <http://sourceforge.net/projects/bamtools/>
- GATK: [http://www.broadinstitute.org/gsa/wiki/index.php/The\\_Genome\\_Analysis\\_Toolkit](http://www.broadinstitute.org/gsa/wiki/index.php/The_Genome_Analysis_Toolkit)

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	NA00001	NA00002
20	14370	rs6054257	G	A	29	0	NS=3;DP=14;AF=0.5;DB;H2	GT:GQ:DP:HQ	0	0:48:1:51
20	17330	.	T	A	3	q10	NS=3;DP=11;AF=0.017	GT:GQ:DP:HQ	0	0:49:3:
20	1110696	rs6040355	A	G,T	67	0	NS=2;DP=10;AF=0.333,0.667;AA=T;DB	GT:GQ:DP:HQ	1	2:21:6:
20	1230237	.	T	.	47	0	NS=3;DP=13;AA=T	GT:GQ:DP:HQ	0	0:54:7:
20	1234567	microsat1	G	D4,IGA	50	0	NS=3;DP=9;AA=G	GT:GQ:DP	0	1:35:4

Variants: VCF

- VCFTools: <http://vcftools.sourceforge.net/>
- VcfCTools: <https://github.com/AlistairNWard/vcfCTools>

# Project timeframe (approximate)

- Phase 1
  - Raw data, alignments available
  - Integrated variant set available
  - Phase 1 analysis paper by end of 2011
- Phase 2
  - Raw data mid-December 2011
  - Read mapping, variant calling early 2012
- Phase 3
  - Samples end March 2012
  - Data Summer 2012
  - Call sets end of 2012, Final paper 2013?
- End of the project

# Fraction of variant sites present in an individual that are NOT already represented in dbSNP

Date	Fraction <u>not</u> in dbSNP
February, 2000	98%
February, 2001	80%
April, 2008	10%
February, 2011	2%
October 2011 (now)	<1%