

# The 1000 Genomes Project

Advanced Information  
Laura Clarke



# Command Line Tools

- Samtools <http://samtools.sourceforge.net/>
- VCFTools <http://vcftools.sourceforge.net/>
- Tabix <http://sourceforge.net/projects/samtools/files/tabix/>
  - (Please note it is best to use the trunk svn code for this as the 0.2.5 release has a bug)
  - svn co <https://samtools.svn.sourceforge.net/svnroot/samtools/trunk/tabix>



# Sequence Data

- Fastq files
  - @ERR050087.1 HS18\_6628:8:1108:8213:186084#2/1
  - GGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGGGTTAGG
  - +
  - DCDHKHKKIJGNNHIJIIKLLMCLKMAILIJH3K>HL1I=>MK.D
  - <http://www.1000genomes.org/faq/what-format-are-your-sequence-files>



# Alignment Data

- BAM files
- ERR052835 163 11 60239 0 100M = 60609 469
- <http://samtools.sourceforge.net/>

NAME	DESCRIPTION
QNAME	Query NAME of the read or read pair
FLAG	Bitwise FLAG (pairing, strand, mate strand etc
RNAME	Reference Sequence NAME
POS	1-Based leftmost POSition of clipped alignment
MAPQ	MAPping Quality (Phred-scaled)
CIGAR	Extended CIGAR string (operations: MIDNSHP)
MRNM	Mate Reference NaMe ('=' if same as RNAME)
MPOS	1-Based leftmost Mate POSition
ISIZE	Inferred Insert SIZE
SEQ	Query SEQUENCE on the same strand as the reference
QUAL	Query QUALity (ASCII-33=Phred base quality)



# Alignment data: Extended Cigar Strings

Cigar has been traditionally used as a compact way to represent a sequence alignment. BAM files contain an extended version of this cigar string

Operations include

**M** - match or mismatch

**I** - insertion

**D** - deletion

SAM extends these to include

**S** - soft clip

**H** - hard clip

**N** - skipped bases

**P** - padding

E.g. Read: ACGCA-TGCAGTtagacgt

Ref: ACTCAGTG----GT

Cigar: 5M1D2M2I2M7S



# Variant Call Data

- VCF Files
- TAB Delimited Text Format

NAME	DESCRIPTION
CHROM	Chromosome name
POS	Position in chromosome
ID	Unique Identifier of variant
REF	Reference Allele
ALT	Alternative Allele
QUAL	Phred scaled quality value
FILTER	Site filter information
INFO	User extensible annotation
FORMAT	Describes the format of the subsequent fields, must always contain Genotype
Individual Genotype Fields	These columns contain the individual genotype data for each individual in the file



# Variant Call Data

- Headers

```
##fileformat=VCFv4.1
```

```
##INFO=<ID=RSQ,Number=1,Type=Float,Description="Genotype imputation  
quality from MaCH/Thunder">
```

```
##INFO=<ID=AC,Number=.,Type=Integer,Description="Alternate Allele Count">
```

```
##INFO=<ID=AN,Number=1,Type=Integer,Description="Total Allele Count">
```

```
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele, ftp://ftp.  
1000genomes.ebi.ac.uk/vol1/ftp/technical/reference/ancestral_alignments/  
README">
```

```
##INFO=<ID=AF,Number=1,Type=Float,Description="Global Allele Frequency  
based on AC/AN">
```

```
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
```

```
##FORMAT=<ID=DS,Number=1,Type=Float,Description="Genotype dosage  
from MaCH/Thunder">
```

```
##FORMAT=<ID=GL,Number=.,Type=Float,Description="Genotype  
Likelihoods">
```

# Variant Call Data

- Example 1000 Genomes Data
- CHROM 4
- POS 42208061
- ID rs186575857
- REF T
- ALT C
- QUAL 100
- FILTER PASS
- INFO AA=T;AN=2184;AC=1;RSQ=0.8138;AF=0.0005;
- FORMAT GT:DS:GL
- GENOTYPE 0|0:0.000:-0.03,-1.19,-5.00













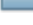







ftp://ftp.1000genomes.ebi.ac.uk

ftp://ftp-trace.ncbi.nih.gov/1000genomes/ftp

Index of ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/

 Up to higher level directory

Name	Size	Last Modified
 CHANGELOG	118 KB	05/01/2012 5/01/2012 12:40:00
 README.alignment_data	12 KB	26/01/2011 26/01/2011 12:00:00
 README.ftp_structure	9 KB	04/04/2011 4/04/2011 12:00:00
 README.pilot_data	3 KB	14/07/2011 14/07/2011 12:00:00
 README.populations	2 KB	18/02/2010 18/02/2010 12:00:00
 README.sequence_data	7 KB	23/07/2011 23/07/2011 19:03:00
 alignment_indices		14/07/2011 14/07/2011 10:53:00
 changelog_details		05/01/2012 05/01/2012 12:40:00
 current.tree	29933 KB	05/01/2012 05/01/2012 12:37:00
 data		04/07/2011 04/07/2011 8:50:00
 phase1		14/07/2011 14/07/2011 14:03:00
 pilot_data		27/07/2011 27/07/2011 12:00:00
 release		12/10/2011 12/10/2011 13:18:00
 sequence.index	27185 KB	20/12/2011 20/12/2011 12:26:00
 sequence_indices		14/11/2011 14/11/2011 10:10:00
 technical		13/12/2011 13/12/2011 10:05:00

Documentation

Raw Data

Phase 1 Data

Pilot Data

Release Data

Technical Data

# Meta Data Formats

- Sequence Index
  - Sequence meta data from ENA
- Alignment Index
  - Location and md5sum for Alignment Files
- BAS
  - Read group level alignment statistics
- HsMetrics
  - Exome alignment statistics based on Picard CalculateHsMetrics



# Sequence Index

- Meta Data File to present information about each fastq file
- Allows easy location of specific subsets of data
- Use to denote specific sequence freezes
- Sequence\_indices directory contains complete history
- Named
  - YYYYMMDD.sequence.index
  - 20120130.sequence.index is most current



Sequence Index	Description	Column	Description
1. Fastq File	Relative path to file	14. Instrument Model	Sequencing Machine Model
2. MD5 checksum	Checksum for file	15. Library Name	
3. Run ID	SRA run id	16. Run Name	
4. Study ID	SRA study id	17. Run Block Name	No Longer used
5. Study Name	SRA study descriptor	18. Insert Size	Estimated Insert Size
6. Center Name	Submission Center	19. Library Layout	Paired or Single ended
7. Submission ID	SRA submission id	20. Paired Fastq	Paired Fastq File
8. Submission Date	Date of Submission	21. Withdrawn	Withdrawn Status
9. Sample ID	SRA Sample ID	22. Withdrawn Date	
10. Sample Name	Coriell Sample name	23. Withdrawn Reason	
11. Population	Population Code	24. Read Count	
12. Experiment ID	SRA Experiment ID	25. Base Count	
13. Instrument Platform	Sequencing Machine Platform	26. Analysis Group	Sequencing Strategy

# Alignment Index

- 6 column file pointing to location of BAM files
- Bam filenames contains majority of information
  - Sample\_name.location.instrument\_platform.alignment\_algorithm.population.analysis\_group.Index\_data.bam
- Alignment index lines contains location and md5 for
  - BAM file
  - BAI file
  - BAS file

# Bas files

- Alignment statistics
- Read group level stats for each alignment
- 21 column file including
  - Read group name
  - Sample name
  - Total Base Count
  - Mapped Base Count
  - Duplicate Base Count



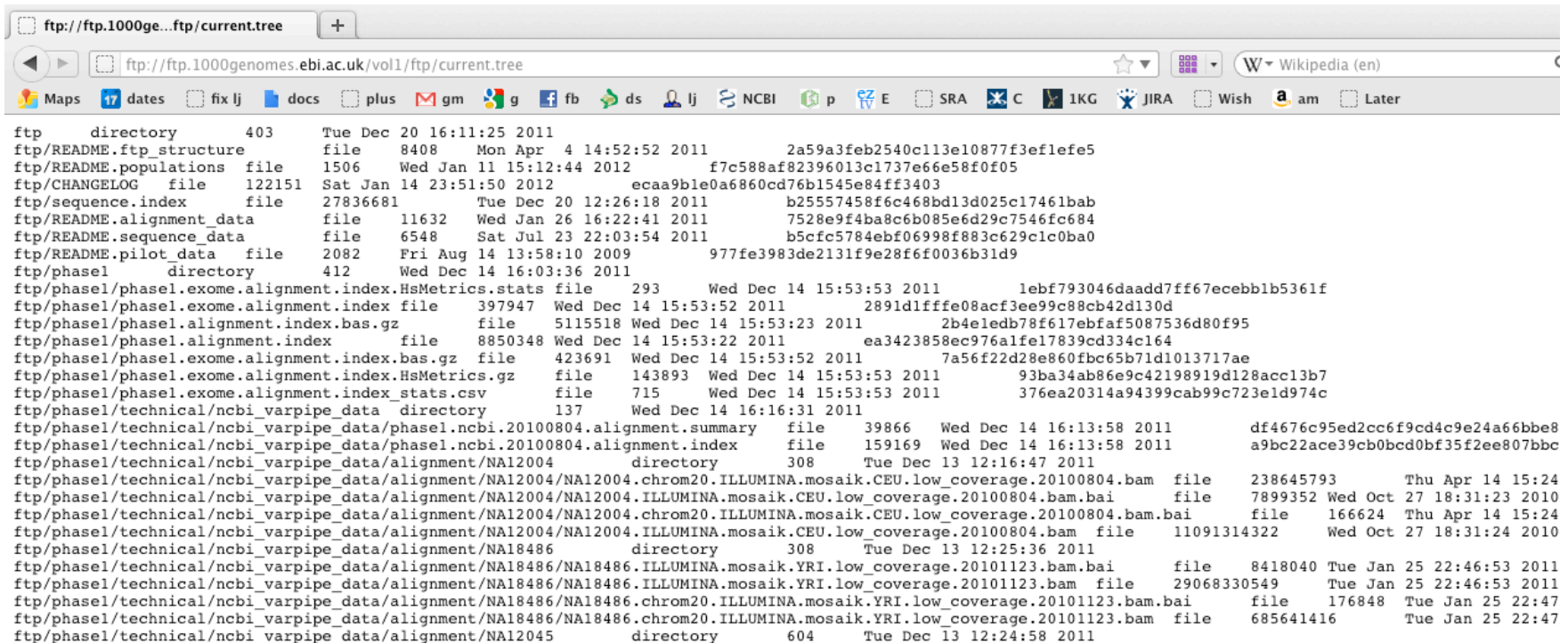
# HsMetrics Files

- Picard Command line tool, CalculateHsMetric
- Used to define completed Exome
- Distributed in gzipped format
- Contains 38 columns like
  - File\_name
  - ON\_BAIT\_BASES
  - MEAN\_BAIT\_COVERAGE
  - PCT\_TARGET\_BASES\_20X



# Finding Data

- Current.tree file
- <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/current.tree>
- Current Tree is updated nightly so can be upto 24 hours out of date



```
ftp directory 403 Tue Dec 20 16:11:25 2011
ftp/README.ftp_structure file 8408 Mon Apr 4 14:52:52 2011 2a59a3feb2540c113e10877f3ef1efe5
ftp/README.populations file 1506 Wed Jan 11 15:12:44 2012 f7c588af82396013c1737e66e58f0f05
ftp/CHANGELOG file 122151 Sat Jan 14 23:51:50 2012 ecaa9b1e0a6860cd76b1545e84ff3403
ftp/sequence.index file 27836681 Tue Dec 20 12:26:18 2011 b2557458f6c468bd13d025c17461bab
ftp/README.alignment_data file 11632 Wed Jan 26 16:22:41 2011 7528e9f4ba8c6b085e6d29c7546fc684
ftp/README.sequence_data file 6548 Sat Jul 23 22:03:54 2011 b5cfc5784ebf06998f883c629c10ba0
ftp/README.pilot_data file 2082 Fri Aug 14 13:58:10 2009 977fe3983de2131f9e28f6f0036b31d9
ftp/phase1 directory 412 Wed Dec 14 16:03:36 2011
ftp/phase1/phase1.exome.alignment.index.HsMetrics.stats file 293 Wed Dec 14 15:53:53 2011 1ebf793046daadd7ff67ecebb1b5361f
ftp/phase1/phase1.exome.alignment.index file 397947 Wed Dec 14 15:53:52 2011 2891d1ffffe08acf3ee99c88cb42d130d
ftp/phase1/phase1.alignment.index.bas.gz file 5115518 Wed Dec 14 15:53:23 2011 2b4e1edb78f617ebfaf5087536d80f95
ftp/phase1/phase1.alignment.index file 8850348 Wed Dec 14 15:53:22 2011 ea3423858ec976a1fe17839cd334c164
ftp/phase1/phase1.exome.alignment.index.bas.gz file 423691 Wed Dec 14 15:53:52 2011 7a56f22d28e860fbc65b71d1013717ae
ftp/phase1/phase1.exome.alignment.index.HsMetrics.gz file 143893 Wed Dec 14 15:53:53 2011 93ba34ab86e9c42198919d128acc13b7
ftp/phase1/phase1.exome.alignment.index_stats.csv file 715 Wed Dec 14 15:53:53 2011 376ea20314a94399cab99c723e1d974c
ftp/phase1/technical/ncbi_varpipe_data directory 137 Wed Dec 14 16:16:31 2011
ftp/phase1/technical/ncbi_varpipe_data/phase1.ncbi.20100804.alignment.summary file 39866 Wed Dec 14 16:13:58 2011 df4676c95ed2cc6f9cd4c9e24a66bbe8
ftp/phase1/technical/ncbi_varpipe_data/phase1.ncbi.20100804.alignment.index file 159169 Wed Dec 14 16:13:58 2011 a9bc22ace39cb0bcd0bf35f2ee807bbc
ftp/phase1/technical/ncbi_varpipe_data/alignment/NA12004 directory 308 Tue Dec 13 12:16:47 2011
ftp/phase1/technical/ncbi_varpipe_data/alignment/NA12004/NA12004.chrom20.ILLUMINA.mosaik.CEU.low_coverage.20100804.bam file 238645793 Thu Apr 14 15:24
ftp/phase1/technical/ncbi_varpipe_data/alignment/NA12004/NA12004.chrom20.ILLUMINA.mosaik.CEU.low_coverage.20100804.bam.bai file 7899352 Wed Oct 27 18:31:23 2010
ftp/phase1/technical/ncbi_varpipe_data/alignment/NA12004/NA12004.chrom20.ILLUMINA.mosaik.CEU.low_coverage.20100804.bam.bai file 166624 Thu Apr 14 15:24
ftp/phase1/technical/ncbi_varpipe_data/alignment/NA12004/NA12004.ILLUMINA.mosaik.CEU.low_coverage.20100804.bam file 11091314322 Wed Oct 27 18:31:24 2010
ftp/phase1/technical/ncbi_varpipe_data/alignment/NA18486 directory 308 Tue Dec 13 12:25:36 2011
ftp/phase1/technical/ncbi_varpipe_data/alignment/NA18486/NA18486.ILLUMINA.mosaik.YRI.low_coverage.20101123.bam.bai file 8418040 Tue Jan 25 22:46:53 2011
ftp/phase1/technical/ncbi_varpipe_data/alignment/NA18486/NA18486.ILLUMINA.mosaik.YRI.low_coverage.20101123.bam file 29068330549 Tue Jan 25 22:46:53 2011
ftp/phase1/technical/ncbi_varpipe_data/alignment/NA18486/NA18486.chrom20.ILLUMINA.mosaik.YRI.low_coverage.20101123.bam.bai file 176848 Tue Jan 25 22:47
ftp/phase1/technical/ncbi_varpipe_data/alignment/NA18486/NA18486.chrom20.ILLUMINA.mosaik.YRI.low_coverage.20101123.bam file 685641416 Tue Jan 25 22:47
ftp/phase1/technical/ncbi_varpipe_data/alignment/NA12045 directory 604 Tue Dec 13 12:24:58 2011
```





# Data Slicing

- All alignment and variant files are indexed so subsections can be downloaded remotely
- Use samtools to get subsections of bam files
  - **samtools view** [http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data/HG01375/alignment/HG01375.mapped.ILLUMINA.bwa.CLM.low\\_coverage.20111114.bam](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data/HG01375/alignment/HG01375.mapped.ILLUMINA.bwa.CLM.low_coverage.20111114.bam) 6:31833200-31834200
- Use tabix to get subsections of vcf files
  - **tabix -h** [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120131\\_omni\\_genotypes\\_and\\_intensities/Omni25\\_genotypes\\_2141\\_samples.b37.vcf.gz](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120131_omni_genotypes_and_intensities/Omni25_genotypes_2141_samples.b37.vcf.gz) 6:31833200-31834200
- You can also use the web Data Slicer interface to do this



# Data Slicing

- VCFtools provides some useful additional functionality on the command line including:
- vcf-compare, comparison and stats about two or more vcf files
- vcf-isec, creates an intersection of two or more vcf files
- vcf-subset, will subset a vcf file only retaining the specified individual columns
- vcf-validator, will validate a particular



# Exercise, Finding Data

1. How many GRCh37 omni vcf files are in technical/working
2. Which exome sample from 20110521 has the highest percentage of targets covered at 20x or greater. You need to look at the 20110521.alignment.index.HsMetrics.gz file to find this
3. Find the exome bam file for this sample
4. Get a slice of this exome bam file between 7:114173990-114175942



# Exercise Answers, Finding Data

```
> wget ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/current.tree
> grep omni current.tree | cut -f1 | grep vcf | grep -v tbi | grep
b37 | wc -l
> 32
> zcat 20110521.exome.alignment.index.HsMetrics.gz | cut
-f1,31 | sort -k2 -n | tail -n1
> HG00737.mapped.illumina.mosaik.PUR.exome.
20110411.bam 0.932651
```



# Exercise Answers, Finding Data

```
>samtools view ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/
phase1/data/HG00737/exome_alignment/
HG00737.mapped.illumina.mosaik.PUR.exome.
20110411.bam 7:114173990-114175942 | tail -n1
> SRR099984.44615561      83      7      114174990      65
76M      =      114174660      -405
GAACCATATTTGGTGTACATAGGCATAAAGAATTTTGCA
TAAAACCCCCTTGTGGGATTTTATTCATACATAGGTT
SD@GIB>BFDDHDCDBBJCAFHHJBBDDEHDBFFDCHJB
<CCC4IIHHIECGCGGGAEeee@AEBH??@H@?CFDBS
RG:Z:SRR099984 NM:i:0 OQ:Z:DE@DEE?
EEBEGEDEGFHHFGHHHHGHHFHHGHHDHHHHHHGHHD
HHGGGHHHHHHHHHHHHHHHHHHHHGHHHHHHHHHH
```

# Command Line Tools



# Variant Effect Predictor

- Predicts Functional Consequences of Variants
- Both Web Front end and API script
- Can provide
  - sift/polyphen/condel consequences
  - Refseq gene names
  - HGVS output
- Can run from a cache as well as Database
- Convert from one input format to another
- Script available for download from:
  - [ftp://ftp.ensembl.org/pub/misc-scripts/Variant\\_effect\\_predictor/](ftp://ftp.ensembl.org/pub/misc-scripts/Variant_effect_predictor/)
  - [http://browser.1000genomes.org/Homo\\_sapiens/UserData/UploadVariations](http://browser.1000genomes.org/Homo_sapiens/UserData/UploadVariations)



# Variant Effect Predictor

- `perl variant_effect_predictor.pl -input 6_381831625_3184704.vcf -sift p -polyphen p -check_existing`
- `less variant_effect_output.txt`

```
#Uploaded_variation Location Allele Gene Feature Feature_type Consequence
cDNA_position CDS_position Protein_position Amino_acids Codons Exi
sting_variation Extra
rs138094825 6:31831667 A ENSG00000204385 ENST00000414427 Transcript
DOWNSTREAM - - - - - rs138094825 -
rs138094825 6:31831667 A ENSG00000204385 ENST00000229729 Transcript
INTRONIC - - - - - rs138094825 -
6_31832657_C/T 6:31832657 T ENSG00000204385 ENST00000229729
Transcript NON_SYNONYMOUS_CODING 1883 1862 621 R/H cGc/cAc -
PolyPhen=possibly_damaging;SIFT=deleterious
```



# Data Slicing

- Use samtools to get subsections of bam files
  - **samtools view** [http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data/HG01375/alignment/HG01375.mapped.ILLUMINA.bwa.CLM.low\\_coverage.20111114.bam](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data/HG01375/alignment/HG01375.mapped.ILLUMINA.bwa.CLM.low_coverage.20111114.bam) 6:31833625-31833704
- Use tabix to get subsections of vcf files
  - **tabix -h** [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120131\\_omni\\_genotypes\\_and\\_intensities/Omni25\\_genotypes\\_2141\\_samples.b37.vcf.gz](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120131_omni_genotypes_and_intensities/Omni25_genotypes_2141_samples.b37.vcf.gz)  
6:31830969-31846823 | **vcf-subset -c HG01375**
- [http://browser.1000genomes.org/Homo\\_sapiens/UserData/SelectSlice](http://browser.1000genomes.org/Homo_sapiens/UserData/SelectSlice)



# Variation Pattern Finder

- Remote or local tabix indexed VCF input
- Discovers patterns of Shared Inheritance
- Variants with functional consequences considered by default
- Web output with CSV and Excel downloads
- [http://browser.1000genomes.org/Homo\\_sapiens/  
UserData/VariationsMapVCF](http://browser.1000genomes.org/Homo_sapiens/UserData/VariationsMapVCF)



# Variation Pattern Finder

- **perl variant\_pattern\_finder.pl** -vcf ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20110521/ALL.chr6.phase1\_integrated\_calls.20101123.snps\_indels\_svsvs.genotypes.vcf.gz -sample\_panel\_file ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20110521/phase1\_integrated\_calls.20101123.ALL.panel -region 6:31830969-31846823 -expand



# Variation Pattern Finder Output

freq	6:31833647_[T]	6:31833660_rs6915800[G]	samples	
freq	ENST00000414427- SPlice_SITE[],ENST0000054 4672- SPlice_SITE[],ENST0000022 9729- SPlice_SITE[],ENST0000037 5562-SPlice_SITE[]	ENST00000414427- NON_SYNONYMOUS_CODING[R/ C],ENST00000229729- NON_SYNONYMOUS_CODING[R/ C],ENST00000544672- NON_SYNONYMOUS_CODING[R/ C],ENST00000375562- NON_SYNONYMOUS_CODING[R/C]	samples	
0.73	REF REF	G A	YRI(3)	NA18933, NA19149, NA19098 and 0 others.
0.27	REF REF	A G	YRI(2)	NA19146, NA19198
0.18	REF REF	A A	LWK(1)	NA19372
0.09	C T	REF REF	CHB(1)	NA18592



# VCF to PED

- LD Visualization tools like Haploview require PED files
- VCF to PED converts VCF to PED
- Will a file divide by individual or population
- [http://browser.1000genomes.org/Homo\\_sapiens/UserData/Haploview](http://browser.1000genomes.org/Homo_sapiens/UserData/Haploview)



# VCF to PED

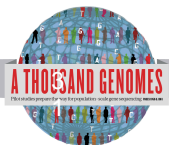
- **perl vcf\_to\_ped\_convert.pl** -vcf [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20110521/ALL.chr6.phase1\\_integrated\\_calls.20101123.snps\\_indels\\_svsvs.genotypes.vcf.gz](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20110521/ALL.chr6.phase1_integrated_calls.20101123.snps_indels_svsvs.genotypes.vcf.gz) -sample\_panel\_file [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20110521/phase1\\_integrated\\_calls.20101123.ALL.panel](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20110521/phase1_integrated_calls.20101123.ALL.panel) -region **6:31830969-31846823** -population **CEU**
- Output should be two files
- 6\_31830969-31846823.info
- 6\_31830969-31846823.ped

# Haploview

- haploview



<http://www.broadinstitute.org/scientific-community/science/programs/medical-and-population-genetics/haploview>



# Access to backend Ensembl databases

- Public MySQL database at
  - `mysql-db.1000genomes.org` port 4272
- Full programmatic access with Ensembl API
  - The 1000 Genomes Pilot uses Ensembl v60 databases and the NCBI36 assembly (this is frozen)
  - The 1000 Genomes main project currently uses Ensembl v63 databases
- <http://jun2011.archive.ensembl.org/info/docs/api/variation/index.html>
- <http://www.ensembl.org/info/docs/api/variation/index.html>
- <http://www.1000genomes.org/node/517>





# Amazon Web Service Cloud

- 1000 Genomes Alignments and Variant files are available in AWS
- AMI image available to run 1000 Genomes Tutorial
- <http://www.1000genomes.org/using-1000-genomes-data-amazon-web-service-cloud>



# Data Availability

- FTP site: <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/>
  - Raw Data Files
- Web site: <http://www.1000genomes.org>
  - Release Announcements
  - Documentation
- Ensembl Style Browser: <http://browser.1000genomes.org>
  - Browse 1000 Genomes variants in Genomic Context
  - Variant Effect Predictor
  - Data Slicer
  - Other Tools



# Exercises, Command Line Tools

5. Get a slice of `HG00737.mapped.illumina.mosaik.PUR.exome.20110411.bam` for `7:114304000-114305000` (FoxP2 exon)
6. Get the equivalent section of the 20110521 release chr 7 genotypes file
7. Use `vcftools vcf-subset` to get the genotypes for HG00737, does HG00737 have any variant sites in this location?
8. Use this piece with tools, the variant effect predictor, the vcf pattern finder
9. Are there any snps with deleterious sift/polyphen consequences?
10. What is the most common pattern of variation in this region?
11. Use the vcf to ped script with `6:31830700-31840700` and population CEU
12. How many different haplotype blocks does the section contain?



# Exercise Answers, Command Line Tools

```
> grep HG00737.mapped.illumina.mosaik.PUR.exome.20110411.bam /  
nfs/1000g-archive/vol1/ftp/current.tree | cut -f1 | grep -v bam. | awk  
'{print "ftp://ftp.1000genomes.ebi.ac.uk/vol1/ "$1}'
```

```
> ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/ftp/phase1/data/HG00737/  
exome_alignment/HG00737.mapped.illumina.mosaik.PUR.exome.  
20110411.bam
```

```
> samtools view ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/ftp/phase1/  
data/HG00737/exome_alignment/  
HG00737.mapped.illumina.mosaik.PUR.exome.20110411.bam
```

```
SRR099984.29321596 163 7 114304108 65 76M = 114304379 346  
GTTTGCTGCAAGGACGATTGTTTATATTTTCACATCGCACTTAATTTCTTGCATCTCTGCCACAAG  
TAGCCAGTT S=??DDBGE@CGGAE@BABIACB?  
A@ACCCGCGCBH=GCGEBAEBCDHHCEIHBBGDHEIHHCGABIAAIHHCGBR RG:Z:SRR099984  
NM:i:0  
OQ:Z:HHHHEHHHHHHHHFHHHHHHHHGEGHGHHHHHHHHHBHFHGFHHHHHHHHFHHHEHHFHHH  
HHHHDBGFGGHHFEHF  
SRR099984.344934 163 7 114304134 59 76M = 114304429 370  
TTTTACATCGCACTTAATTTCTTGCATCTCTGCCACAAGGAGCCAGTTAGGAATTTTTTTTCAATA  
CATTTTCT S>??D?C??B>A6BBB?C>A AFF1ECCB9FECBDAD=CAEG&AGDDBGAB@GFB@@@?  
>B781<=<?@87>=55>5S RG:Z:SRR099984 NM:i:1 OQ:Z:FHEHHHHGGEGFD6EFGGEBDEGG/  
I@EG8GGDBBBE=FBFE,BEEDEDEFFEE@FB@F8:3>@?@D==A?77@77
```



# Exercise Answers, Command Line Tools

```
> tabix -h ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/
20110521/
ALL.chr7.phase1_release_v3.20101123.snps_indels_svsvs.g
notypes.vcf.gz 7:114304000-114305000 > 20110521.vcf
> tabix -h ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/
20110521/
ALL.chr7.phase1_release_v3.20101123.snps_indels_svsvs.g
notypes.vcf.gz 7:114304000-114305000 | vcf-subset -c
HG00737 > HG00737.vcf
> vcf-stats HG00737.vcf
'hom_RR_count' => 16,
```



# Exercise Answers, Command Line Tools

```
>perl variant_effect_predictor.pl -input ~/20110521.vcf -sift p  
-polyphen p --force_overwrite
```

```
> grep SIFT variant_effect_output.txt
```

```
> rs182138317 7:114304331 A  
ENSG00000128573 ENST00000393489 Transcript  
NON_SYNONYMOUS_CODING 1949 1567 523 A/T  
Gcc/Acc -
```

PolyPhen=possibly\_damaging;SIFT=deleterious

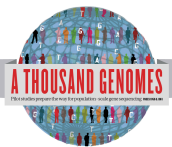


# Exercise Answers, Command Line Tools

```
> perl variant_pattern_finder.pl -vcf ~/20110521.vcf -sample  
ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20110521/  
phase1_integrated_calls.20101123.ALL.panel -region  
7:114304000-114305000
```

This produces a tsv file which can be view in a spreadsheet program

```
7:114304563_rs1378771[C] 7:114304630_rs1378772[A] 7:114304969_rs2396765[T]  
> 14.38 - - C|T - A|T - - - - - T|C TSI(30)
```



# Exercise Answers: Command Line Tools

```
> perl vcf_to_ped_convert.pl -vcf 20110521.vcf -sample  
ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20110521/  
phase1_integrated_calls.20101123.ALL.panel -region  
6:31830700-31840700 -population CEU
```

```
> ls ./
```

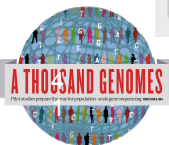
```
> 6_31830700-31840700.info
```

```
> 6_31830700-31840700.ped
```

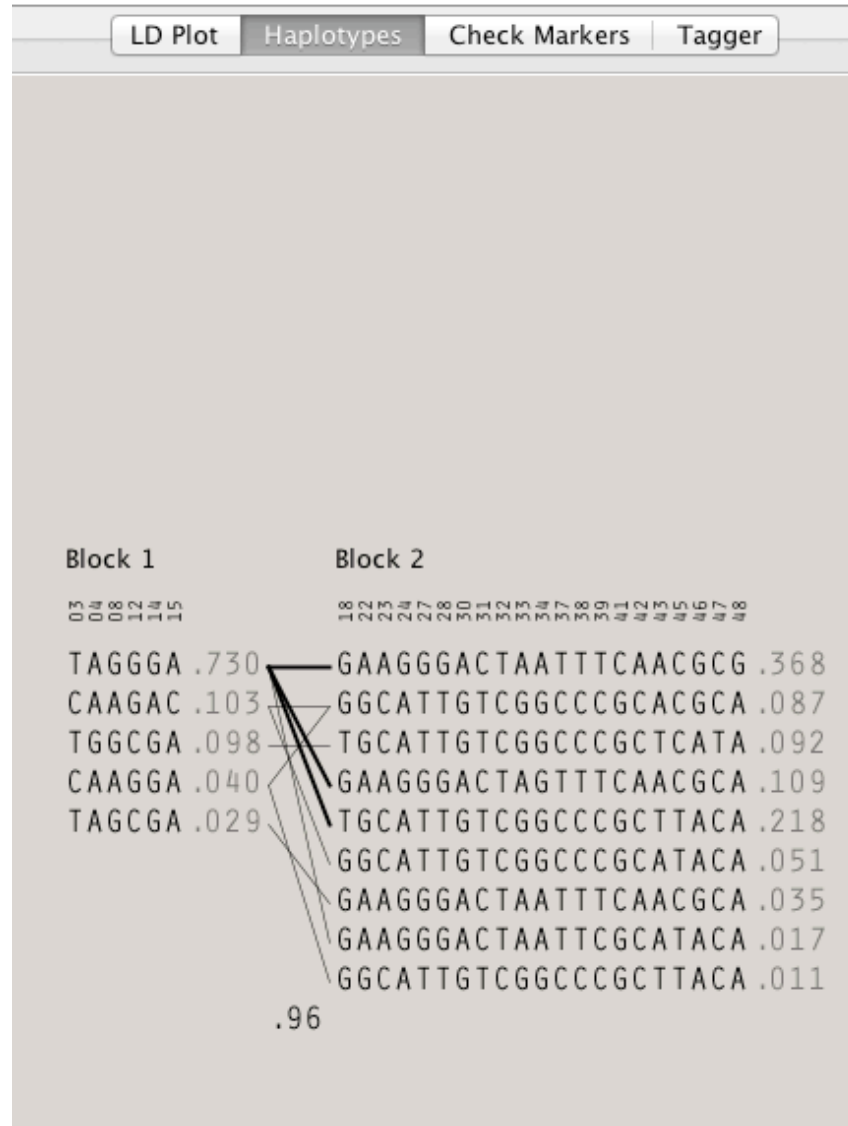




# Exercise Answers



# Exercise Answers



# Questions

Please send any future questions about this presentation and any other material on our website to [info@1000genomes.org](mailto:info@1000genomes.org)



<http://www.1000genomes.org/using-1000-genomes-data>



# Thanks

- The 1000 Genomes Project Consortium
- Paul Flicek
- Richard Smith
- Holly Zheng Bradley
- Ian Streeter
- David Richardson

