# Generation Gap: How existing bioinformatics resources are adapting to high-throughput sequencing
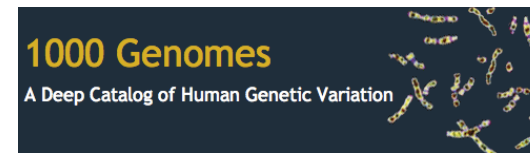
Paul Flicek

Vertebrate Genomics

EMBL-EBI

# Further Evolution of Large-scale Genome Sequencing

- 2000: Human genome working drafts
- Data unit of approximately 10x coverage of human
  - 10 years and cost about $3 billion

- 2008: Major genome centers can sequence the same number of base pairs **every 4 days**
  - 1000 Genome project launched
  - World-wide capacity dramatically increasing

- 2009: **Every 4 hours ($25,000)**
- 2010: **Every 14 minutes ($5,000)**
  - Illumina HiSeq2000 machine produces 200 gigabases per 8 day run (BGI have 128)





1000 Genomes
A Deep Catalog of Human Genetic Variation



illumına

# Large-scale genome sequencing

- Today
  - 1000 Genomes, Cancer Genomes, exomes
  - Personal Genomes, Celebrity Genomes, Family Genomes
  - Others
- Soon
  - Thousands of cancer genomes
  - UK 10K
  - Diagnostic laboratories
  - Much, much more
- Results
  - Astronomical amounts of data
  - Catalogs of human variation and mutation

# How is next generation sequencing data impacting major bioinformatics resources

- We have always attached a diverse community of users
  - From absolute beginners to ninjas
  - All need support
- Sequencing data is opening up new experiments and driving the transition to human as the model organisms
  - Variation data is the largest component of this change
- Multiple challenges
  - Data access for those who what big and small pieces
  - Annotation and management of the resulting discoveries
  - Your genome is unique and so is everyone else's genome*

\* Identical twins not included

EMBL-EBI

# 1000 Genomes Project: Primary goals

- Overall: Create a deep catalogue of human variation to provide a better baseline to underpin human genetics

- Discover shared variation (shared = not private to individual) and characterise by allele frequency
  - Aim for effectively all (not just a lot of) common variation
    - For example: any variant down to 1% minor allele frequency in a population in the accessible genome has a 95% chance of being identified
    - The pilot project and simulations will help to determine the precision of this statement
  - Structural variants as well as SNPs
    - Accessible because the project will used paired-end sequencing reads
  - Deeper discovery in gene regions, down to 0.5% to 0.1% MAF

EMBL-EBI

# 1000 Genomes Project: Outcomes

- A public database of essentially all SNPs and detectable CNVs with allele frequency >1% in each of multiple human population samples
- Pioneer and evaluate methods for:
  - Generating data from next-generation sequencing platforms
  - Exchanging and combining data and analytical methods
  - Discovering and genotyping SNPs and CNVs data
  - Imputation with and from next generation sequencing data
- Produce an open resource building on HGP, HapMap etc.
  - A control set for sequencing disease samples

- All data publicly available, cell lines available
  - Anonymised samples without phenotypes

EMBL-EBI

# 1000 Genomes Project Design and Progress

- Three pilot projects
    - Deep sequence two trios
    - Low coverage (~2X) 60 individuals from each of three populations (180 individuals total)
    - Gene capture for 1000 genes in about 700 individuals

- Pilot data collected in 2008; analysis now finished; paper now submitted to "a major journal"

- Full project data collection and analysis underway

EMBL-EBI

# Pilot Project SNP Discovery



- 84% of novel SNPs to a single population
  - 4% in all populations
- FDR: <5% for SNPs and <10% for small indels

# Genome-wide SNP distributions



- High variation rate
  - HLA region
  - Telomeres
- Low variation in other places
  - notably 3p21

EMBL-EBI

# Variation around genes



- Heterozygosity is lowest in middle exons
- Diversity is proportional to divergence
  - Functional constraint is the driving force of gene diversity

EMBL-EBI

# Value of additional samples for variant discovery comparing exon and low coverage



- 220 LC individuals to find 99% of synonymous variants
- 320 LC individuals for 98.5% of non-synonymous

EMBL-EBI

Pilot project 180 samples
Extension to 1,100 samples summer 2010
1900 samples end 2010, 2500 samples end 2011

Major population groups comprised of subpopulations of ~100 each

# 1000 Genomes data by populations

| Population | Sequence (gigabases) | Total Coverage |
|---|---:|---:|
| ASW | 645 | 215x |
| CEU | 2368 | 789x |
| CHB | 1135 | 378x |
| CHS | 168 | 56x |
| GBR | 141 | 47x |
| JPT | 1841 | 614x |
| LWK | 1087 | 362x |
| MXL | 216 | 72x |
| TSI | 1257 | 419x |
| YRI | 1534 | 511x |

Total number of base pairs as of 11 June – 10.4 TB (12.5 TB including pilot projects)
Approximately 3500x total genome coverage

EMBL-EBI

# Putting this scale of data into perspective

- Size of EMBL/Genbank in April 2008 at the start of the 1000 Genomes Project: 235,135,312,328 nucleotides

- The 1000 Genomes project routinely produces the equivalent amount of sequence *every three days*
  - This is only a fraction of world-wide sequence capacity

- Data sizes in biology are now on the same order as those common in physics and astronomy

EMBL-EBI

# 1KG Data storage infrastructure

Today

Circa April 2008

# Challenges

- ## Access
  - Data size, storage and transfer
  - Providing access to other researchers that want to use the data

- ## Annotation of variant data
  - Incorporating published and curated information
  - Integrating data that is collected on the genome index

- ## Most human research data cannot be openly released
  - How much diagnostic data should be released?
  - From research to clinical practice

EMBL-EBI

# 1000 Genomes Project: Data Flow



- Developed organically with many loops to relatively smooth system that takes data from sequencing machine to FTP site in about 1 month

EMBL-EBI

# The 1000 Genomes data infrastructure

- Most aspects are running relatively smoothly
    - The pilot project produced about 100,000 sequence and other data files (there are now hundreds of thousands more)
    - Reseqtrack knows where the file is, what has been done to it, potential problems, related result files
- Accurate data transfer and bandwidth remain significant problems
    - File corruption during transfer is still relatively common
    - EBI bandwidth demands have increased about four fold over the course of the project
- The groups using this data are still mostly those within the 1000 Genomes project
    - The demand is growing beyond the project participants

EMBL-EBI

# ReseqTrack System: Pipeline overview



Input and event details

ReseqTrack

Store Location Of Output Files

run_event_pipeline.pl

LSF

runner.pl

event

event program

STORAGE

**source**forge  FIND AND DEVELOP OPEN SOURCE SOFTWARE

Find Software   Develop   Create Project   Blog   Site Support   About

SourceForge.net > Develop > ReseqTrack

ReseqTrack

Summary | Files | Support | **Develop** | Tracker | Mailing Lists | Forums | Code

Code

Programming Languages: Perl

Repositories

SVN browse code, statistics, last commit on 2010-06-18

```
svn co https://reseqtrack.svn.sourceforge.net/svnroot/reseqtrack
reseqtrack
```

L. Clarke

EMBL-EBI

# www.1000genomes.org

Project information with regular news updates

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

# ftp.1000genomes.ebi.ac.uk

30-50 terabytes of data

Mostly in data formats that have just been invented
and almost no one has heard of or knows how to use

EMBL-EBI

# 1000 Genomes

## A Deep Catalog of Human Genetic Variation

Home

### Search 1000Genomes

[          ] Go

e.g. gene BRCA2 or AL032821.2.1.143563

### 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.

**SeqAlignView** →
Shows read-depth data alongside SNPs

Other sites using Ensembl software...

### Press Release

### December 2008

**Browser displays SNP calls on CEU and YRI high coverage individuals from Pilot2**

- View sample data
- EBI Mirror
- NCBI Mirror

### The 1000 Genomes Browser

**Ensembl-based browser provides early access to 1000genomes data**

In order to facilitate immediate analysis of the 1000genomes data by the whole scientific community, this browser (based on Ensembl) integrates the SNP calls and read coverage from this December 2008 release. All of this data has been submitted to dbSNP, and once rsid's have been allocated, will be absorbed into the UCSC and Ensembl browsers according to their respective release cycles. Until that point **any SNP id's on this site are temporary and will NOT be maintained.**

### Links

**1000 Genomes** →
More information about the 1000 Genomes Project on the 1000 genomes main site.

**1000 Genomes Wiki** →
Browse the 1000 Genomes Wiki.

The 1000 Genomes Project is an international collaborative project described at 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

1000 Genomes release 3 - May 2010 © EBI   e!mpowered

About 1000Genomes | Contact Us | Help

1000 Genomes Browser Home Page
http://browser.1000genomes.org

EMBL-EBI

# 1000 Genomes
## A Deep Catalog of Human Genetic Variation

Home > Human

Location: 6:133,042,209-133,101,683

**Location-based displays**
- Whole genome
- Chromosome summary
- Region overview
- **Region in detail**
- Comparative Genomics
  - Genomic alignments (0
  - Synteny (0)
- Genetic Variation
  - Resequencing (6)
  - Linkage Data
- Markers

- Configure this page
- Add custom data to page
- Export data
- Bookmark this page

**Chromosome 6: 133,042,2**

Assembly excepti... chromosome 6

Assembly excepti...

« Region overview

Location: 6 13304220

**Resembl settings**

Chromosome bands
Ensembl/Havana g...

Reference
Consensus
Contigs
Ensembl/Havana g...
< VNN1
Known p
CCDS set
1KG:CEU:PILOT1
1KG:CEU:PILOT3
1KG:CHB+JPT:PILOT1
1KG:CHB:PILOT3
1KG:CHD:PILOT3
1KG:JPT:PILOT3
1KG:LWK:PILOT3
1KG:TSI:PILOT3
1KG:YRI:PILOT1
1KG:YRI:PILOT2
1KG:YRI:PILOT3
%GC

Gene Legend
Variation Legend

Reverse s
Known p
Downstr
Non-syn
Splice si
There are cu
Ensembl Hon

Reads

Coverage 60

Reference
Consensus
Sequences

Reference only displayed for less than 0.2 Kb.
Consensus only displayed for less than 0.2 Kb.
Sequences is displayed only for regions less then 0.2 Kb (314)

ℹ **Configuring the display**

You currently have the overview panel and 70 tracks on the main panel turned off. To change the tracks you are displaying, use the "**Configure this page**" link on the left.

cation
ntext
nes

P tracks
o
pulation

22

About 1000Genomes | Contact Us | Help

EMBL-EBI

# Expanding data availability with the cloud

- Amazon Web Services
  - The final 1000 Genomes Pilot alignment files (BAMs) are now loaded into the Amazon EC2 cloud and have been formally announced last week
  - Some files had to be split to accommodate the 5 Gb max file size of the S3 storage
- Anyone can use the data with standard AWS costs per computer hour (as low as 8.5¢ per CPU hour)
- We will be developing publicly accessible applications within the cloud environment and expect that others will as well

Stein *Genome Biology* 2010, 11:207
http://genomebiology.com/2010/11/5/207

Genome **Biology**

**REVIEW**

The case for cloud computing in genome informatics

Lincoln D Stein*

Software
**Searching for SNPs with cloud computing**
Ben Langmead[*†], Michael C Schatz[†], Jimmy Lin[‡], Mihai Pop[†] and Steven L Salzberg[†]

**Open Access**

Addresses: *Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, 615 North Wolfe Street, Baltimore, Maryland 21205, USA. †Center for Bioinformatics and Computational Biology, University of Maryland, College Park, MD 20742, USA. ‡The iSchool, College of Information Studies, University of Maryland, College Park, MD 20742, USA.

Correspondence: Ben Langmead. Email: blangmea@jhsph.edu

EMBL-EBI

# AWS Public Data Sets



We have just launched a complete Ensembl genome browser mirror within EC2 (http://useast.ensembl.org).

EMBL-EBI

# Challenges

- Access
  - Data size, storage and transfer
  - Providing access to other researchers that want to use the data

- **Annotation of variant data**
  - **Incorporating published and curated information**
  - **Integrating data that is collected on the genome index**

- Most human research data cannot be openly released
  - How much diagnostic data should be released?
  - From research to clinical practice

EMBL-EBI

# Ensembl

- Ensembl's mission is to enable genomic science by providing high-quality, integrated annotation on vertebrate genomes within a consistent and accessible infrastructure.

- Creating and providing core value-added data sets
    - High-quality evidence-based gene sets
    - Multiple alignments
    - Gene homology and paralogy relationships
    - Genome variation including SNPs, genotypes and CNV/SV data
    - Integrative analysis of genome regulation

- Roadmap includes extensive support for data on multiple individuals
    - Human cell lines, mouse strains
    - Favouring integrated information

EMBL-EBI

# Phenotype annotation - Genomic

- Genome wide association study data on 672 phenotypes
- Currently over 60,000 phenotype annotations
  - All high-quality, curated and publication based
- Data is growing with every Ensembl release

# Annotation of the variation catalog in Ensembl

- Incorporated variation annotations representing 134 distinct phenotypes

- Reference 186 publications

- Currently 1120 variations with annotated information
  - All high-quality, curated and publication based

- Data is growing with every Ensembl release

genome.gov
National Human Genome Research Institute
National Institutes of Health

Home | About NHGRI | Newsroom | Staff

Research | Grants | Health | Policy & Ethics | Educational Resources | Careers & Training

Home › About NHGRI › About the Office of the Director › Office of Population Genomics › OPG: A Catalog of Published Genome-Wide Association Studies

**Office of Population Genomics**

Overview ; A Catalog of Genome-Wide Association Studies ; Research Programs ; Recent Publications ; Meetings & Workshops ; Notices & Funding Opportunities ; Contact

A Catalog of Published Genome-Wide Association Studies

Potential etiologic and functional implications of genome-wide association loci for human diseases and traits
Click here to read our recent *Proceedings of the Academy of Sciences (PNAS)* article on catalog methods and analysis.

Go to the Catalog

Published Genome-Wide Associations through 3/2009, 398 published GWA at p ≤ 5 x 10⁻⁸

The genome-wide association study (GWAS) publications listed here include only those att nucleotide polymorphisms (SNPs) in the initial stage. Publications are organized from most to online publication if available. Studies focusing only on candidate genes are excluded from this PubMed literature searches, daily NIH-distributed compilations of news and media reports, and database of GWAS literature (HuGE Navigator).

SNP-trait associations listed here are limited to those with p-values < 1.0 x 10⁻⁵. Note that we associations meeting this p-value threshhold. Multipliers of powers of 10 in p-values are rounde allele frequencies are rounded to two decimals. Standard errors are converted to 95 percent co frequencies, p-values, and odds ratios derived from the largest sample size, typically a combin recorded below if reported; otherwise statistics from the initial study sample are recorded. Od to OR > 1 for the alternate allele. Where results from multiple genetic models are available, we prioritized effect sizes (OR's or beta-coefficients) as follows: 1) genotypic model, per-allele estimate; 2) genotypic model, heterozygote estimate, 3) allelic model, allelic estimate.

Published Genome-Wide Associations (view) ›
Credit: Darryl Leja and Teri Manolio

Gene regions corresponding to SNPs were identified from the UCSC Genome Browser. Gene names are those reported by the authors in the original paper. Only one SNP within a gene or region of high linkage disequilibrium is recorded unless there was evidence of independent association.

Occasionally the term "pending" is used to denote one or more studies that we identified as an eligible GWAS, but for which SNP information has not yet been extracted; studies of CNVs are also noted as pending.

PNAS

Potential etiologic and functional implications of genome-wide association loci for human diseases and traits

Lucia A. Hindorff[a,1], Praveen Sethupathy[b,1], Heather A. Junkins[a], Erin M. Ramos[a], Jayashri P. Mehta[c], Francis S. Collins[b,2], and Teri A. Manolio[a,2]
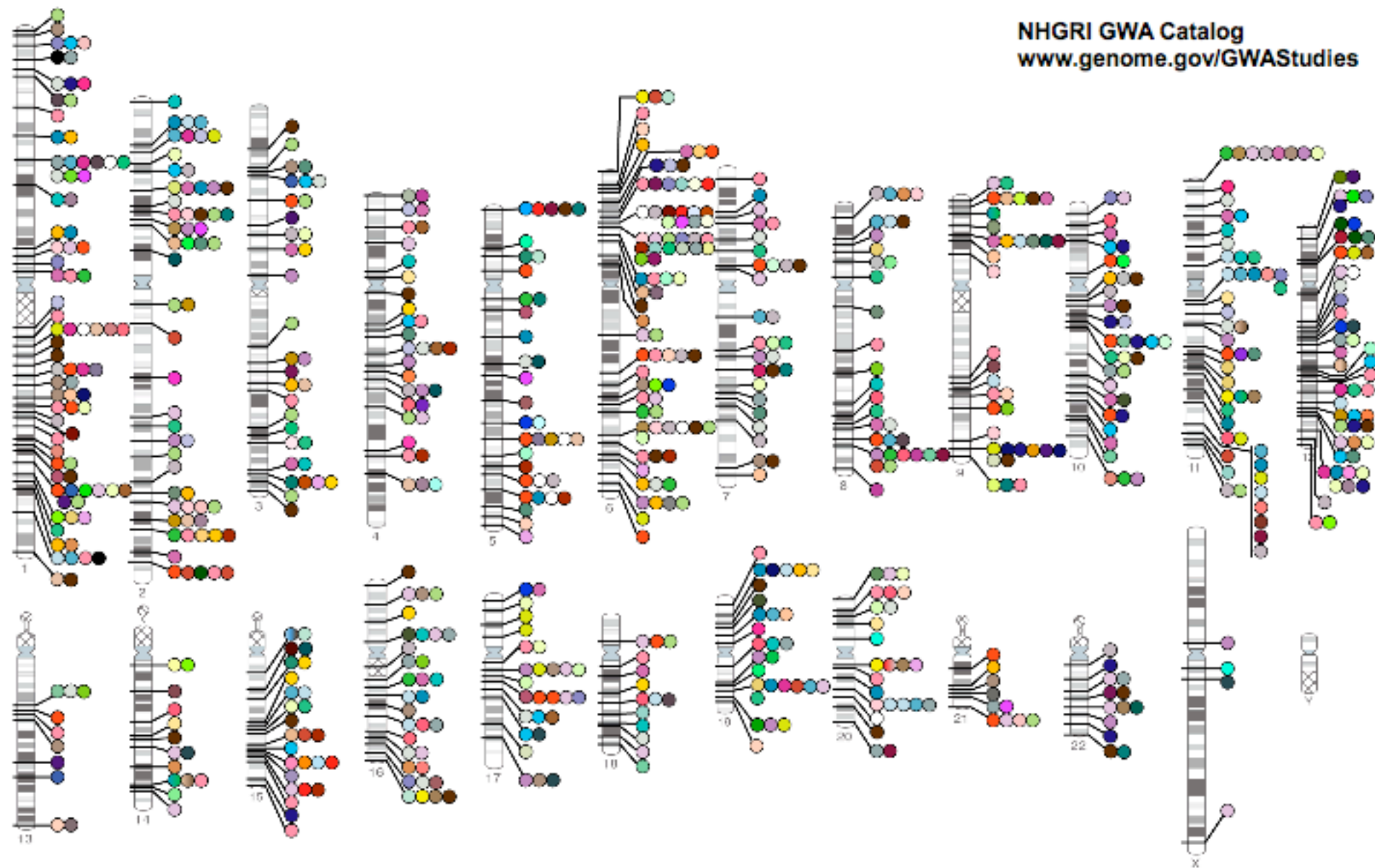
[a]Office of Population Genomics, [b]Genome Technology Branch, National Human Genome Research Institute, and [c]National Center for Biotechnology Information, National Institutes of Health, Bethesda, MD 20892

EMBL-EBI

# Published Genome-Wide Associations through 3/2010, 779 published GWA at $p \leq 5\times10^{-8}$ for 148 traits

Genome

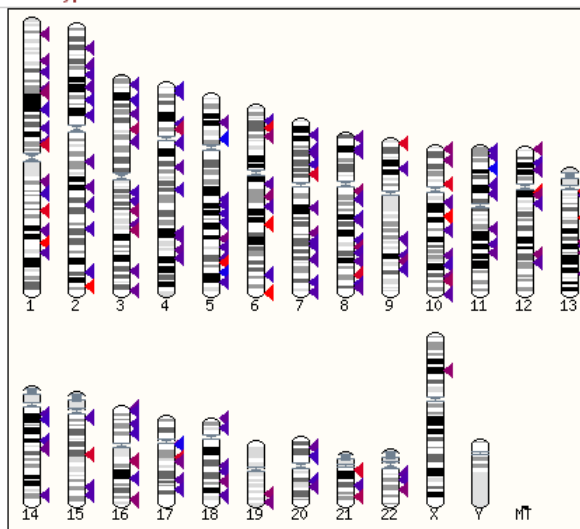**Location-based displays**
- **Whole genome**
- Chromosome summary
- Region overview
- Region in detail
- Comparative Genomics
  - Alignments (image) (51)
  - Alignments (text) (51)
  - Multi-species view (47)
  - Synteny (14)
- Genetic Variation
  - Resequencing (2)
  - Linkage Data
- Markers
- Other genome browsers
  - UCSC
  - NCBI

🔧 Configure this page

📇 Manage your data

📄 Export data

⭐ Bookmark this page

## Karyotype

Whole genome  *help*

**Location of variants associated with phenotype Crohn's Disease:**



Click on the image above to jump to a chromosome, or click and drag to select a region

**Colour Scale:**



| 1.0 | 3.0 | 4.0 | 5.0 | 7.0 | 9.0 | >10.0 |
|---|---|---|---|---|---|---|
| (Least Significant P Value) | | | | | | (Most Significant P Value) |

**Feature Information:**

| Genomic location(strand) | Name(s) | Located in gene(s) | Associated Gene(s) | Associated Phenotype(s) | P value (negative log) |
|---|---|---|---|---|---|
| 1:15435156-15445156(1) | rs6659639 | ENSG00000189337 (RP1-21O18.1) | BC036877 | Crohn's Disease | 4.3 |
| 1:39159153-39169153(1) | rs10493084 | | POU3F1 | Crohn's Disease | 3.9 |
| 1:49458801-49468801(1) | rs3118223 | ENSG00000186094 (AGBL4) | FLJ11588 | Crohn's Disease | 3.1 |
| 1:64226541-64236541(1) | rs2819130 | | | Crohn's Disease | 3.7 |
| 1:67365292-67375292(1) | rs1925411 | ENSG00000152763 (WDR78) | | Crohn's Disease | 5.4 |
| 1:67367061-67377061(1) | rs1983860 | ENSG00000152763 (WDR78) | | Crohn's Disease | 5.4 |
| 1:67433535-67443535(1) | rs12131222 | ENSG00000198160 (MIER1) | | Crohn's Disease | 1.5 |

Done

FoxyProxy: Patterns

EMBL-EBI

# Phenotype annotation - Gene-based

- Over 1400 LSDBs on the Human Genome Variation Society website
    - Data integration with central resources has been challenging

- Locus Reference Genomic Sequences
    - An informatics solution
        - Stable, community-determined sequence
        - Collaboration with the NCBI & Gen2phen
    - Extension and generalisation of NCBI's RefSeqGene project
    - http://www.lrg-sequence.org

Dalgleish, et al. *Genome Medicine* 2010, 2:24

EMBL-EBI

Genome **Medicine**

**CORRESPONDENCE**    **Open Access**

# Locus Reference Genomic sequences: an improved basis for describing hur

Raymond Dalgleish[1]*, Paul Flicek[2], Fiona Cunningham[2]
William M McLaren[2], Pontus Larsson[2], Brendan W Vaug
Peter EM Taschner[7], Johan T den Dunnen[7], Andrew Dev

**EDITORIAL**

nature genetics

## Conventional wisdom

Recent agreement on stable reference sequences for reporting human genetic variants now allows us to mandate the use of the allele naming conventions developed by the Human Genome Variation Society.

**B**y agreement between stakeholders and two principal databases, it has been proposed (R. Dalgleish *et al., Genome Med.* **2**, 24, 2010, doi:10.1186/gm145) that human genetic variants be reported relative to a new set of stable reference sequences, "Locus Reference, Genomic" (LRG, pronounced "large" http://www.lrg-sequence.org/page.php). These sequences have been developed from the initial NCBI RefSeqGene concept and are provided by NCBI and EBI according to agreed rules and in consultation with community users of locus-specific genetic information and locus-specific databases. It is anticipated that the LRG will be stable and supported for many years, long enough to serve as a bridge between existing and future clinical gene tests.

age, resequencing and marker association studies and so keep allele descriptions commensurate with the method by which their data were generated.

The LRG reference sequences should be used in conjunction with standard HGNC gene abbreviations (http://www.genenames.org/) that we already require as a condition of publication. All human genetic variants must now be described—in abstracts and at first use—in accordance with the Human Genome Variation Society (HGVS) conventions (http://www.hgvs.org/mutnomen/) also as a condition of publication. We continue to encourage authors to use HGVS nomenclature for unambiguous reference in all tables and figures and throughout the

EMBL-EBI

**Locus·Reference·Genomic**  [Find LRGs]  Home  Contact  Contributions  FAQ

## Download
You can download the LRG specification document here.

Get the latest version of the LRG XML schema definition and XSLT stylesheets to output HTML and text from the FTP site.

## Contribute
- help@lrg-sequence.org for help and support on the technical issues concerning LRGs (e.g. the XML schema) and the LRG website
- request@lrg-sequence.org for requests to create new LRGs
- feedback@lrg-sequence.org for feedback on the LRG specification

**LRG**

LRG sequences provide a stab framework for reporting mutation ID and core content that ne

We would encourage you to get i convert your RefSeqGene

View a list of LR

To date, there is no internationally recognized reference-sequence standard for reporting sequence varia the NCBI and EBI, as part of the GEN2PHEN consortium, are collaborating with the community of research LSDB curators, mutation consortia etc., to define stable genomic reference sequences called "Locus Refere "LRG". A foundation for this effort is NCBI's RefSeqGene project.

A LRG will provide a stable genomic DNA sequence for a region of the human genome. This sequence ne exactly to a known allele of a gene, but can be idealised to provide a practical working framework. The s will never change, so the unique identifier will not be versioned. The annotation of each LRG is separated particularly to represent exons and coding regions of standard RNA products and their translations as ap updatable section for other biological information such as alternative transcripts, location on the current genome, etc. In particular, the fixed section contains a stable identifier, the genomic, cDNA and amino a as coordinates for the transcript, exon, start and stop codons. The updatable section contains chromoso mapping information for the LRG as well as genomic annotation, database cross-references and alternati amino-acid numbering systems.

EBI and NCBI are committed to developing the technical solutions, as well as computational and visual to sequences. This will enable all the information reported on an LRG to be integrated with the human geno sequence.

For more information on the specification, see the LRG publication:

**Locus Reference Genomic sequences: an improved basis for describing human DNA variants**, Dalgle Med. 2010, 2:24 [View]

or refer to the specification document

This website will list existing LRG sequences and has a FTP site for downloading LRGs. If you would like please email us at feedback@lrg-sequence.org. To create an LRG for your region of interest, please cont sequence.org.

EMBL-EBI    NCBI    GEN2PH

© LRG 2009

---

**Locus·Reference·Genomic**  [Find LRGs]  Home  Contact  Contributions  FAQ

## LRG Search Results
14 results

**LRG_2**
Corresponding HGNC gene symbol: COL1A2
Last updated on: 2010-02-12
View corresponding genomic location: (GRCh37) 7:94018873-94062544 [Ensembl] [NCBI] [UCSC]

**LRG_8**
Corresponding HGNC gene symbol: SCN1A
Last updated on: 2010-02-12
View corresponding genomic location: (GRCh37) 2:166843670-166935149 [Ensembl] [NCBI] [UCSC]

**LRG_13**
Corresponding HGNC gene symbol: CALCA
Last updated on: 2010-02-22
View corresponding genomic location: (GRCh37) 11:14986215-14998832 [Ensembl] [NCBI] [UCSC]

**LRG_6**
Corresponding HGNC gene symbol: ATP1A2
Last updated on: 2010-02-12
View corresponding genomic location: (GRCh37) 1:160080548-160115381 [Ensembl] [NCBI] [UCSC]

**LRG_12**
Corresponding HGNC gene symbol: FKBP10
Last updated on: 2010-03-16
View corresponding genomic location: (GRCh37) 17:39963962-39981469 [Ensembl] [NCBI] [UCSC]

**EMBL-EBI**

# LRGs in Ensembl

# Integrating live external data sources

- SNPedia
  - Wiki-based system for editing information about SNP annotations
    - Current data on 12418 SNPs
  - Licensed under a Creative Commons Attribution-Noncommercial-Share Alike license
  - www.snpedia.com
- Realtime updates in Ensembl



EMBL-EBI

# SNP Effect Prediction tool

- Calculates the effect of SNPs in the context of Ensembl genes and regulatory features
  - Web and API interface
  - Code back-ported to support NCBI36 assembly
  - Programmatic support for tab-delimited and VCF files

- Previously SNP effects were pre-computed for all Ensembl species with variation databases and known SNPs
  - Supports all species and arbitrary SNPs
  - Easily integrated into analysis pipelines

McLaren, et al. *Bioinformatics.* 2010

EMBL-EBI

McLaren, et al. *Bioinformatics.* 2010

# Challenges

- Access
  - Data size, storage and transfer
  - Providing access to other researchers that want to use the data

- Annotation of variant data
  - Incorporating published and curated information
  - Integrating data that is collected on the genome index

- **Most human research data cannot be openly released**
  - How much diagnostic data should be released?
  - From research to clinical practice

EMBL-EBI

# The European Genome-phenome Archive

- Secure storage and authorised access to all types of data sets that might be generated in the context of research into molecular medicine
  - Sequence; Genotypes
  - Transcriptomics; Proteomics
  - Phenotype data
- Enable the collection of larger cohorts and maximisation of resource use
  - Sequencing capacity is increasing dramatically
  - Analysis capacity is increasing more slowly

# EGA Data Acceptance and Access

- Access decisions will remain with the data generating body
  - Distributed model
  - Transparency to the data generators
  - EGA manages the access granted
  - Users can also be restricted to particular collections within a study
- EGA is the European peer database to dbGAP (NCBI)
  - dbGAP has adopted a more centralised model of data access decisions
  - We plan data exchange of meta data and more extensive discussions are on going to increase data discoverability
  - Working toward a common application for both databases to lower administrative burden

# Community Benefits of the EGA

- Data subject to access controls is a burden and it limits the number of researchers that will reuse the resources
  - This may slow the pace of science and prevent serendipitous discovery
- However…
  - Five years ago accessing this type of data was impossible
  - Now it is just incredibly difficult
  - This is real progress
  - Complicated or overly onerous data access agreements are more likely to be ignored

OPINION

## The delay in sharing research data is costing lives

Josh Sommer

It is not uncommon for potentially life-saving research data to be published years after being generated. But the setback to progress caused by the delay in releasing data is troublesome for people who selflessly participate in trials and desperately await new therapies. Scientists need to feel greater urgency to share their findings quickly, and they need additional avenues to facilitate this process.

"Making science work fast enough for patients will require researchers to treat information with greater urgency. Surely, if anyone knew that he or she possessed life-saving data, he or she would act swiftly to share it, just as an intelligence officer would rush to report evidence of an impending terrorist attack."

Nature Medicine, July 2010

EMBL-EBI

# EGA Consortium Page for WTCCC

# Study Page for WTCCC T2D

# WTCCC T1D Data Access Page

# Beyond research toward medical practice

- Needs:
  - Consistent, traceable data generation and analysis routines
  - Robust annotation based on public information sources such as those at the EBI
  - Reporting into medical records

- Data storage:
  - Probably not necessary for primary data as costs drop
  - Individual variant catalogs are already much smaller than MRI data
  - May prevent some liability issues

EMBL-EBI

# Enabling clinical services

- Multiple commercial clinical services built on annotation/Ensembl
  - Alamut - Mutation Interpretation Software - http://www.interactive-biosoftware.com/

# Enabling clinical services

- Multiple commercial clinical services built on annotation/Ensembl
  - BENCH - array CGH platform - http://www.cartagenia.com/

# LRGs are already part of HL7

## HL7 VERSION 2 IMPLEMENTATION GUIDE: CLINICAL GENOMICS; FULLY LOINC-QUALIFIED GENETIC VARIATION MODEL, RELEASE 1 (1ST INFORMATIVE BALLOT)

ORU^R01

HL7 Version 2.5.1

APRIL, 2009

| Chapter Chair: | Amnon Shabo<br>IBM |
|---|---|
| Chapter Chair and Contributing Author: | Mollie Ullman-Cullere<br>Partners HealthCare Center for Personalized Genetic Medicine and Partners Healthcare |
| Chapter Chair: | Phil Pochon<br>Covance |
| Project Chair and Principal Author: | Stan Huff<br>Intermountain Healthcare |
| Project Chair and Contributing Author: | Grant Wood<br>Intermountain Healthcare |
| Contributing Author | Clement McDonald<br>Lister Hill Center for Biomedical Communication, National Library of Medicine |
| Contributing Author | Yan Heras |

### 6.1.5 Reference Sequences (required)

Reference sequences are the baseline from which variation is reported. For example, sequence variants are identified in a patient by comparing the patient's DNA sequence to a reference sequence standard, used in the laboratory. Typically, differences between the patient and reference sequence are called sequence variation and are cataloged, interpreted and reported. Documentation of the reference sequence used is becoming increasingly important for normalization of results between laboratories. To meet this need NCBI is cataloging reference sequences used in clinical testing in the Core Nucleotide Database and can be referred to through the RefSeq identifiers. In collaboration with NCBI, the European BioInformatics Institute (EBI) is also developing a database of reference sequences called Locus Reference Genomic Sequences (LRG). The standard is still in draft status. Importantly, NCBI's RefSeq and EBI's LRG will contain the same reference sequences, annotations and cross references to each other.

### 6.1.6 RefSeq

| TABLE 6-3 - REFSEQ | |
|---|---|
| Code sets, vocabularies, terminologies and nomenclatures that need to be constrained: | RefSeq |
| Minimum attributes of the component: | RefSeq ID |
| Other Comments: | National Center for Biotechnology Information (NCBI) Reference Sequences contained in Core Nucleotide database. Available at: http://www.ncbi.nlm.nih.gov/sites/entrez?db=nuccore. Accessed: March 6, 2008. |

| TABLE 6-3 - LRG | |
|---|---|
| Code sets, vocabularies, terminologies and nomenclatures that need to be constrained: | LRG |
| Minimum attributes of the component: | LRG ID |
| Other Comments: | Locus Reference Genomic Sequences an emerging standard led by the European Bioinformatics Institute |

EMBL-EBI

Annotating the variation catalogue created by the 1000 Genomes projects (and other similar projects) will be one of the major future challenges in human genomics

The results of this annotation will change the way that medicine is practised. And will impact society.

EMBL-EBI

# Acknowledgements

1000 Genomes
A Deep Catalog of Human Genetic Variation

http://www.1000genomes.org

EMBL-EBI