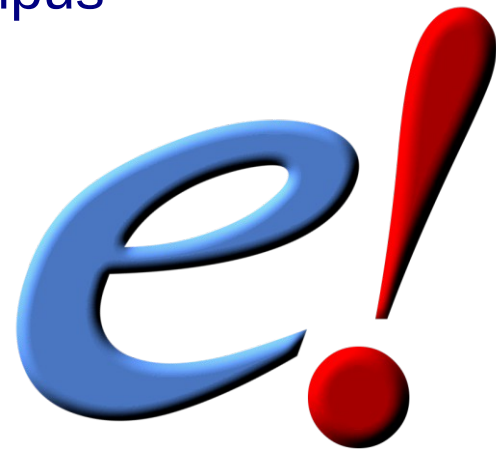# Ensembl

Javier Herrero

Vertebrate Genomics Team

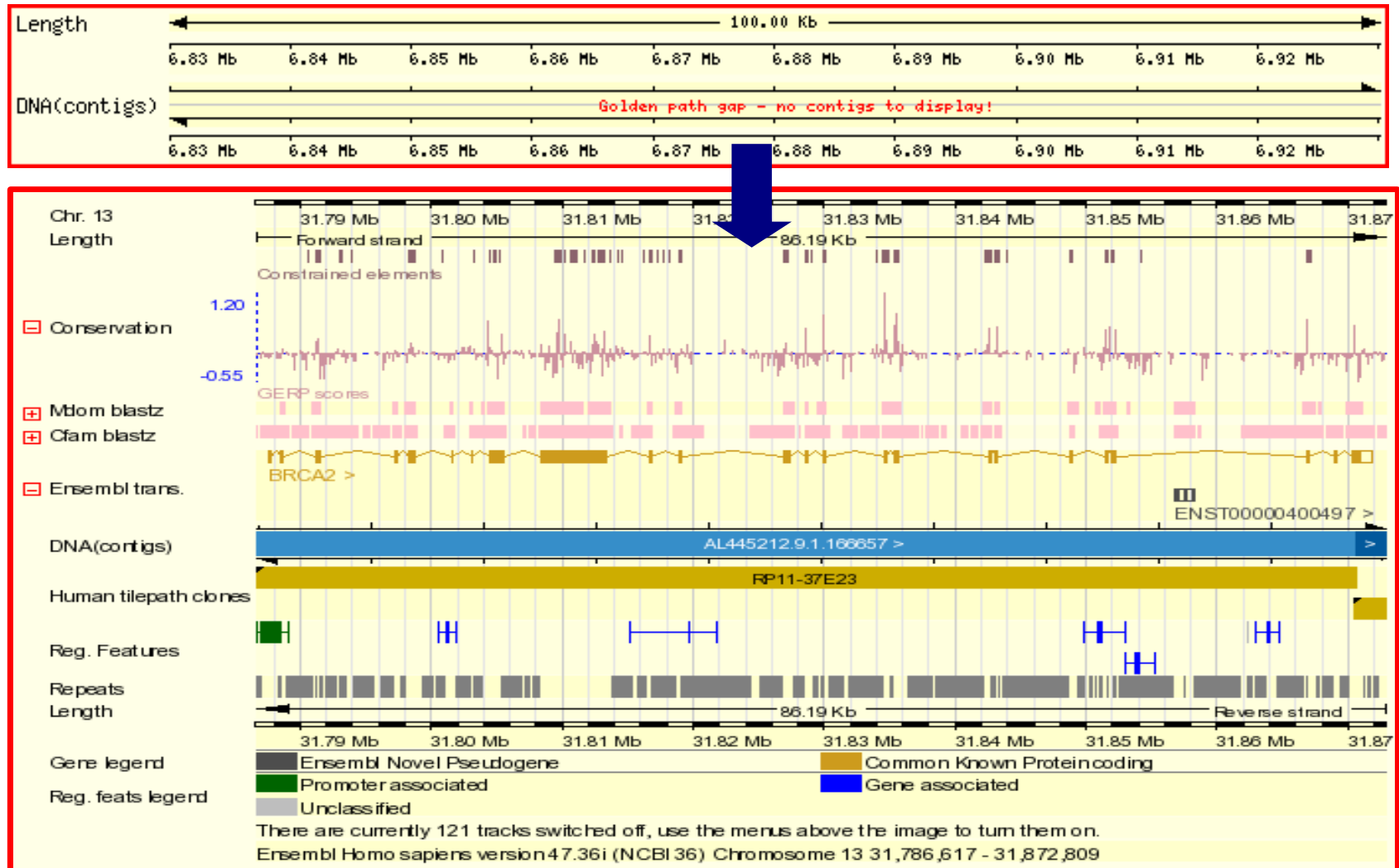EMBL-EBI

Wellcome Trust Genome Campus

Hinxton, UK

# Ensembl mission

To enable genomic science by providing high-quality, integrated annotation on vertebrate genomes within a consistent and accessible infrastructure.

# Ensembl Concept

- Collaborative project of the European Bioinformatics Institute (EBI) and the Wellcome Trust Sanger Institute

- Provides annotation and analysis of chordate genomes

- Open by design
    - Code is BSD, not GNU
    - All data is freely available

- Continuously developed and comprehensively updated 5 times a year

- Diverse skills across the project

- Technology adopted and used by many other projects
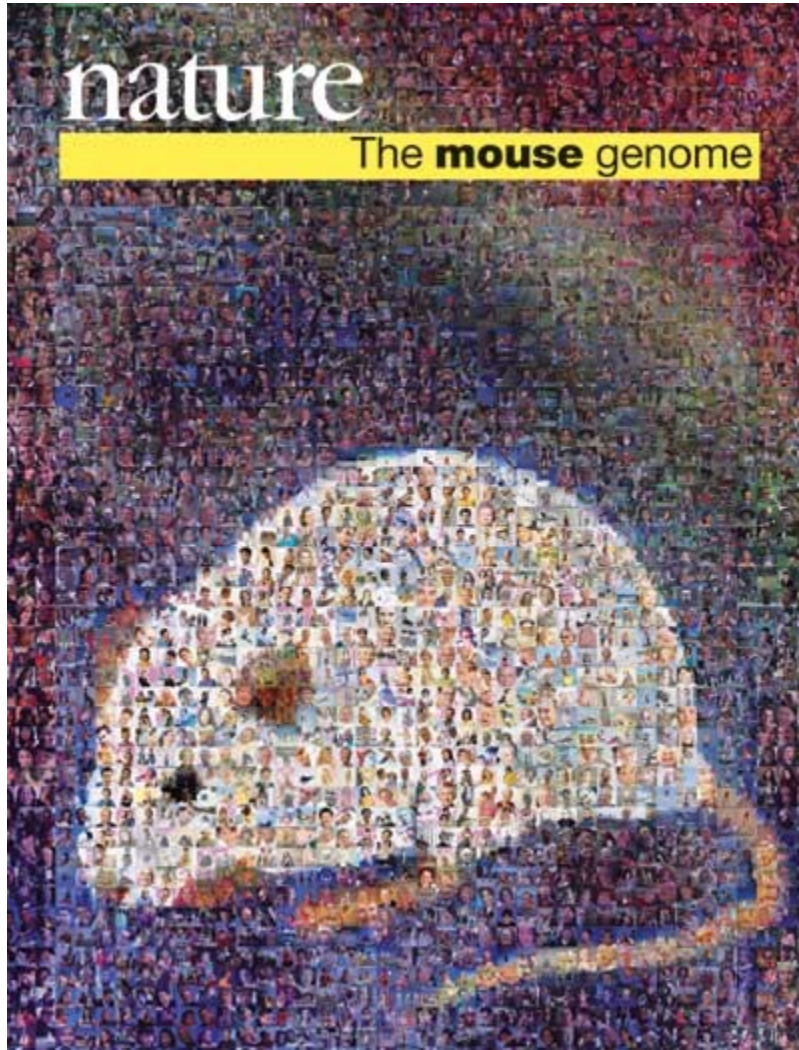
# Human genome (Feb 2001)



Nature, 15th Feb 2001



Science, 15th Feb 2001

# Mouse genome (Dec 2002)



- 2[nd] mammal genome

- model organism in lab

- 14% smaller than human

- 40% can be aligned to human

- 5% under purifying selection

- 0.5 substitutions per site, twice as many in the mouse lineage

wellcome trust **sanger** institute

e!

EMBL-EBI

# Rat genome (Apr 2004)



- 3$^{rd}$ mammal, 2$^{nd}$ rodent

- Similar number of genes in all 3 species

- 40% eutherian specific seq.

- 30% rodent specific seq., mostly repeats

- At least half of unaligned seq. is rat-specific repeats

Nature 428, 493-521
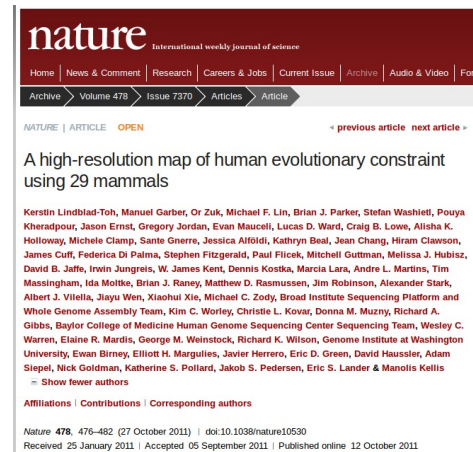
# More and more genomes



2004



2005



2007



2008



2011

More than 50 vertebrate genomes have been "fully" sequenced

Laurasatheria

Afrotheria

Xenarthra

Marsupials

Birds

Fish

Glires

Primates

outspecies

Dasypus_novemcinctus
Choloepus_hoffmanni
Echinops_telfairi
Procavia_capensis
Loxodonta_africana
Sorex_araneus
Erinaceus_europaeus
Pteropus_vampyrus
Myotis_lucifugus
Canis_familiaris
Mustela_putorius_furo
Ailuropoda_melanoleuca
Felis_catus
Equus_caballus
Sus_scrofa
Bos_taurus
Tursiops_truncatus
Vicugna_pacos
Ochotona_princeps
Oryctolagus_cuniculus
Ictidomys_tridecemlineatus
Cavia_porcellus
Dipodomys_ordii
Rattus_norvegicus
Mus_musculus
Tupaia_belangeri
Otolemur_garnettii
Microcebus_murinus
Tarsius_syrichta
Callithrix_jacchus
Macaca_mulatta
Nomascus_leucogenys
Pongo_abelii
Gorilla_gorilla
Pan_troglodytes
Homo_sapiens

Monodelphis_domestica
Sarcophilus_harrisii
Macropus_eugenii
Ornithorhynchus_anatinus
Taeniopygia_guttata
Gallus_gallus
Meleagris_gallopavo
Anolis_carolinensis
Pelodiscus_sinensis
Xenopus_tropicalis
Latimeria_chalumnae
Oreochromis_niloticus
Tetraodon_nigroviridis
Takifugu_rubripes
Xiphophorus_maculatus
Oryzias_latipes
Gasterosteus_aculeatus
Gadus_morhua
Danio_rerio
Petromyzon_marinus
Ciona_savignyi
Ciona_intestinalis
Drosophila_melanogaster
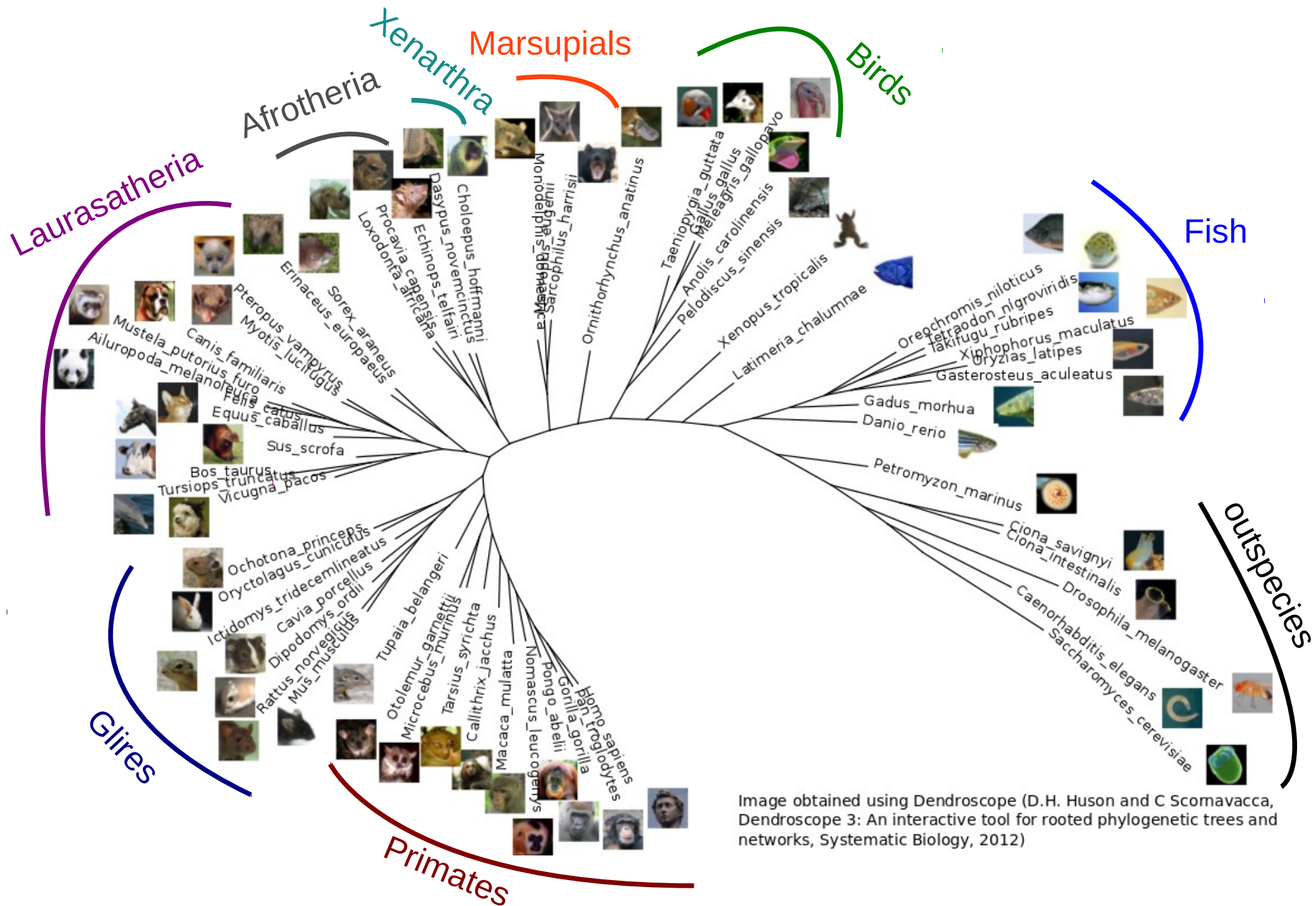Caenorhabditis_elegans
Saccharomyces_cerevisiae

Image obtained using Dendroscope (D.H. Huson and C Scornavacca, Dendroscope 3: An interactive tool for rooted phylogenetic trees and networks, Systematic Biology, 2012)

wellcome trust
sanger institute

e!

EMBL-EBI

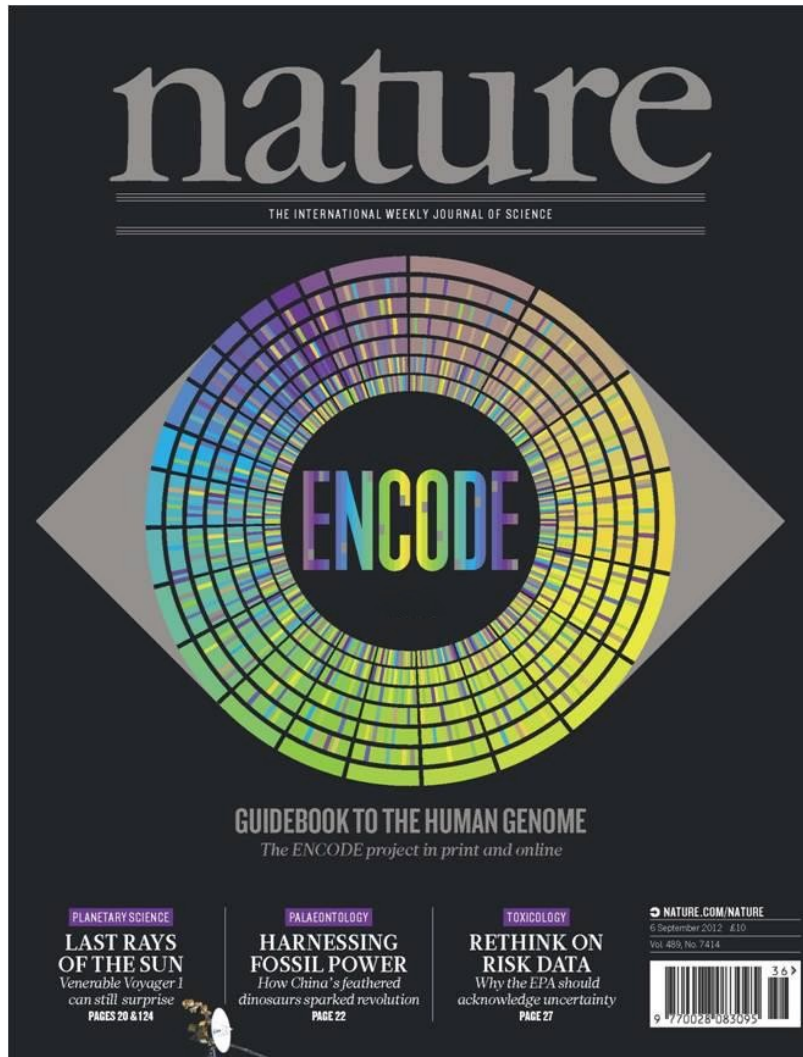# 1000 genomes pilot project (Oct 2010)



3 pilots:

- – 179 ind. low-cov
- – 2 trios high-cov
- – 679 ind. exon only

Covers >95% of variants of any individual

Each individual carries 250-300 loss-of-function variants; 50-100 implicated in inherited disorders

# ENCODE project: ENCyclopedia Of DNA Elements



Nature 489, 57–74 (06 Sept 2012)

- 1640 data sets
- 147 different cell types
- 400+ authors

- Activity in 80% of the genome
- Evidence of negative selection (in aggregate) in primate-specific elements
- Could classify the genome into 7 different chromatin states
- Transcription ↔ Histone + TF

# Ensembl: What do you get?

Genome Annotation

- – Protein coding gene structure
    - • Consistent with genome, predicted across all vertebrates
    - • Manual annotations (human, mouse, zebrafish, MHC)
- – RNA genes (including miRNA)
    - • Consistent with genome, predicted in across mammals
- – Additional identifiers per genes (Xref)
    - • Affymetrix, EntrezGene, Uniprot…

Variation, Comparative & Functional Genomics

- – Genome alignments
    - • Blastz, Blat, Pecan (multiple alignments), EPO
    - • Homologues between genomes
    - • Protein trees
- – Variants (SNPs), CNVs, strains, genotypes
- – ChIP-chip, ChIP-seq, segmentations

Infrastructure

- – Website, Data mining tool, database and data dumps
- – Portable, extendible, open source system with database, DAS, API, website, pipeline

# Ensembl release cycle



genebuilder (8) ~ 3 months

core (4)

regulation (4)

variation (5)

compara (6)

mart (2)

web (5)

release!!

release coordinator + assistant

# Ensembl groups

- **GeneBuilders**: sequence masking, gene building

- **Core**: database schema, stable id mapping

- **Compara**: protein homology, genomic sequence alignments

- **Variation**: SNPs, CNVs, personal/strain genomes

- **Regulation**: probe mapping, functional data, segmentation

- **Web**: web site, new views for new data

- **Outreach**: help, workshops, tutorials

- **Production**: BioMart, coordination

- More people: Research, e! genomes, Zebrafish, Systems...

# Acknowledgements

EMBL

European Commission
Framework Programme 7

EMBL-EBI