



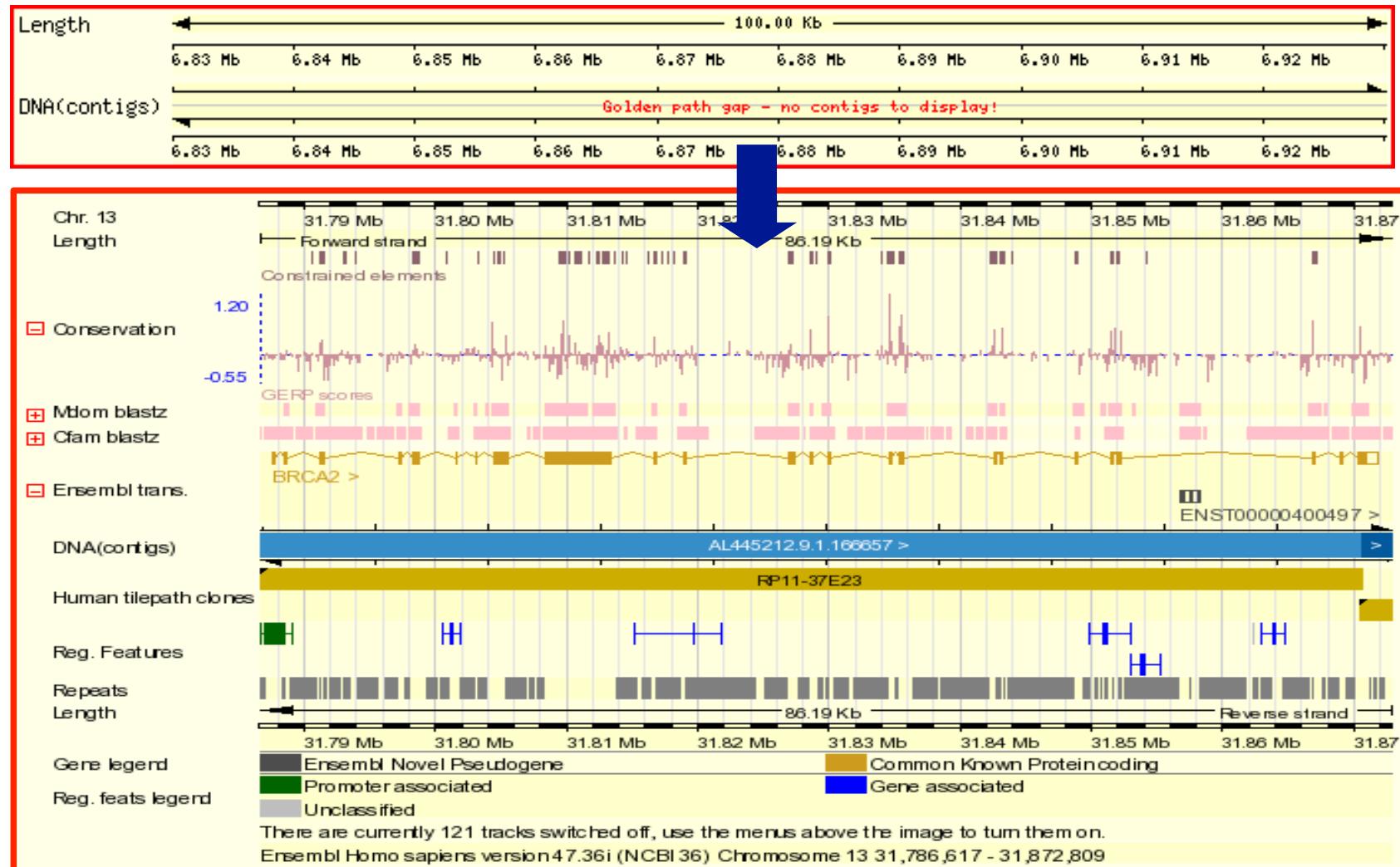
Ensembl

Javier Herrero (*Former Ensembl member*)
Comparative Genomics Project Leader
The Genome Analysis Center (TGAC)
Norwich Research Park
Norwich, UK



Ensembl mission

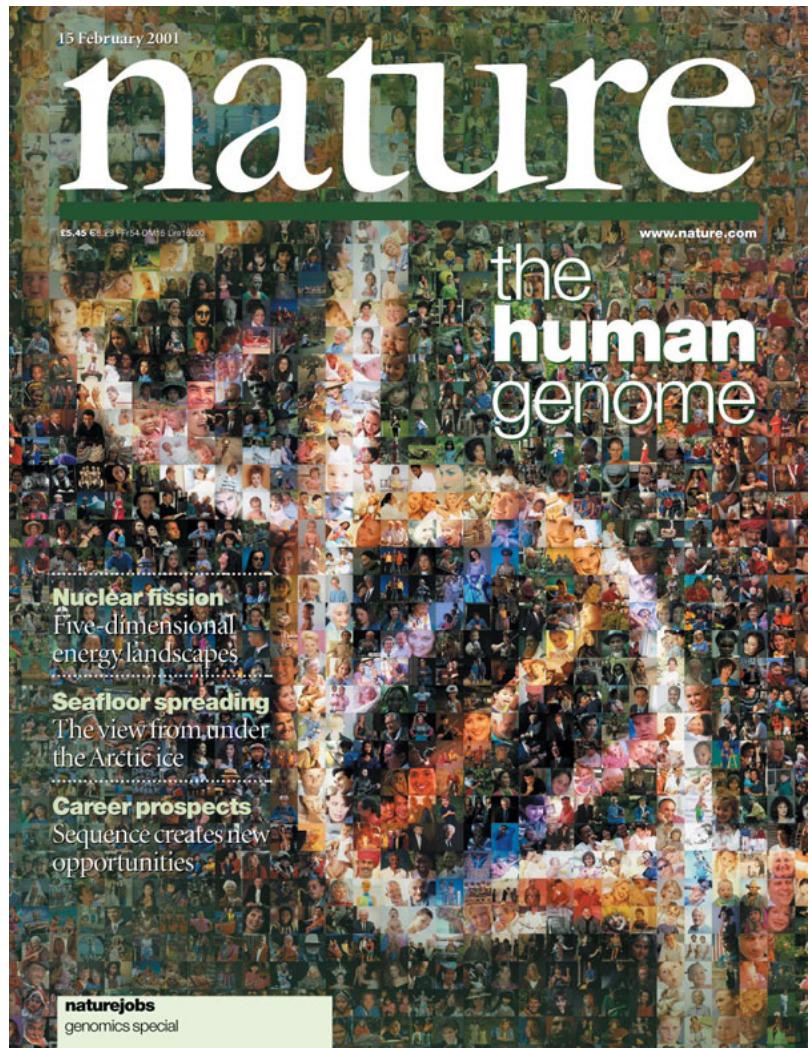
To enable genomic science by providing high-quality, integrated annotation on vertebrate genomes within a consistent and accessible infrastructure.



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

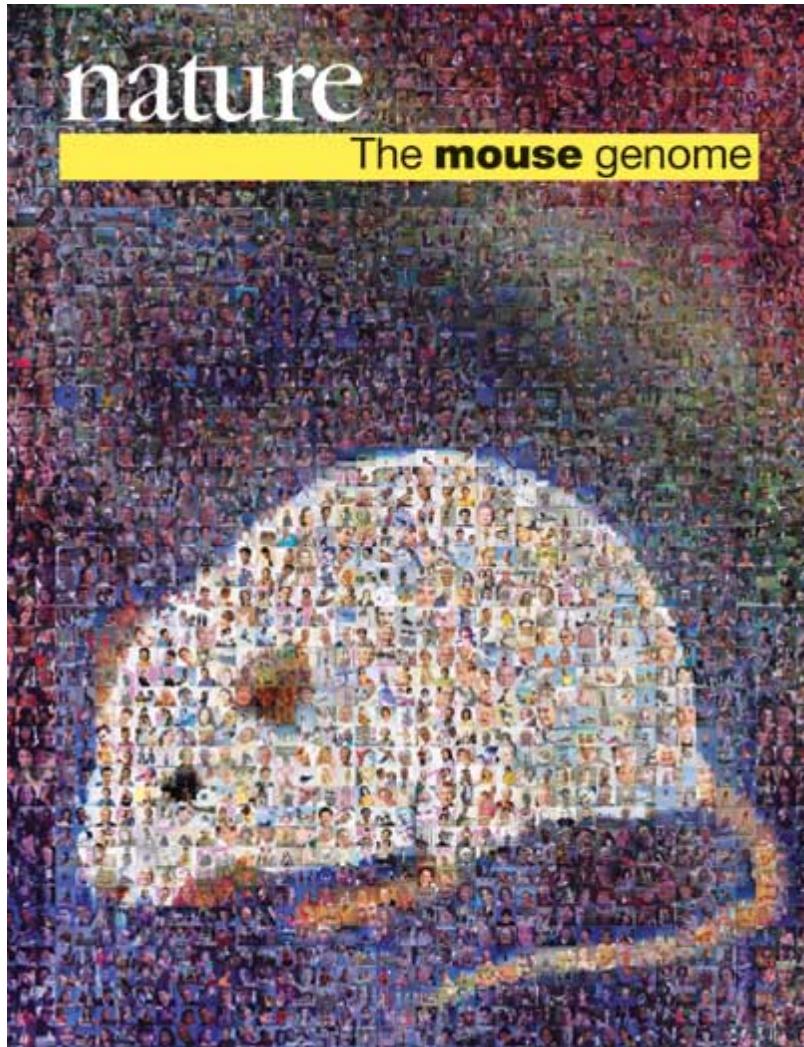
Human genome (Feb 2001)



Nature, 15th Feb 2001

Science, 15th Feb
2001

Mouse genome (Dec 2002)



- 2nd 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

Nature 420, 520-562

Rat genome (Apr 2004)



- 3rd mammal, 2nd 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

A high-resolution map of human evolutionary constraint using 29 mammals

Kersstin Lindblad-Toh, Manuel Garber, Or Zuk, Michael F. Lin, Brian J. Parker, Stefan Washietl, Pouya Kheradpour, Jason Ernst, Gregory Jordan, Evan Mauceli, Lucas D. Ward, Craig B. Lowe, Alisha K. Holloway, Michele Clamp, Sante Gnerre, Jessica Alföldi, Kathryn Beal, Jean Chang, Hiram Clawson, James Cuff, Federica Di Palma, Stephen Fitzgerald, Paul Flicek, Mitchell Gutman, Melissa J. Hubisz, David B. Jaffe, Inren Jungreis, W. James Kent, Dennis Kostka, Marcia Lora, Andre L. Martins, Tim Massingham, Ida Moltó, Brian J. Rooney, Matthew D. Rasmussen, Jim Robinson, Alexander Stark, Albert J. Vitelli, Jiayu Wen, Xiaohui Xie, Michael C. Zody, Broad Institute Sequencing Platform and Whole Genome Assembly Team, Kim C. Worley, Christie L. Kovar, Donna M. Muzny, Richard A. Gibbs, Baylor College of Medicine Human Genome Sequencing Center Sequencing Team, Wesley C. Warren, Elaine R. Mardis, George M. Weinstock, Richard K. Wilson, Genome Institute at Washington University, Ewan Birney, Elliott H. Margulies, Javier Moreno, Eric D. Green, David Haussler, Adam Siepel, Nick Goldman, Katherine S. Pollard, Jakob S. Pedersen, Eric S. Lander & Manolis Kellis

More than 50 vertebrate genomes have been “fully” sequenced

2011

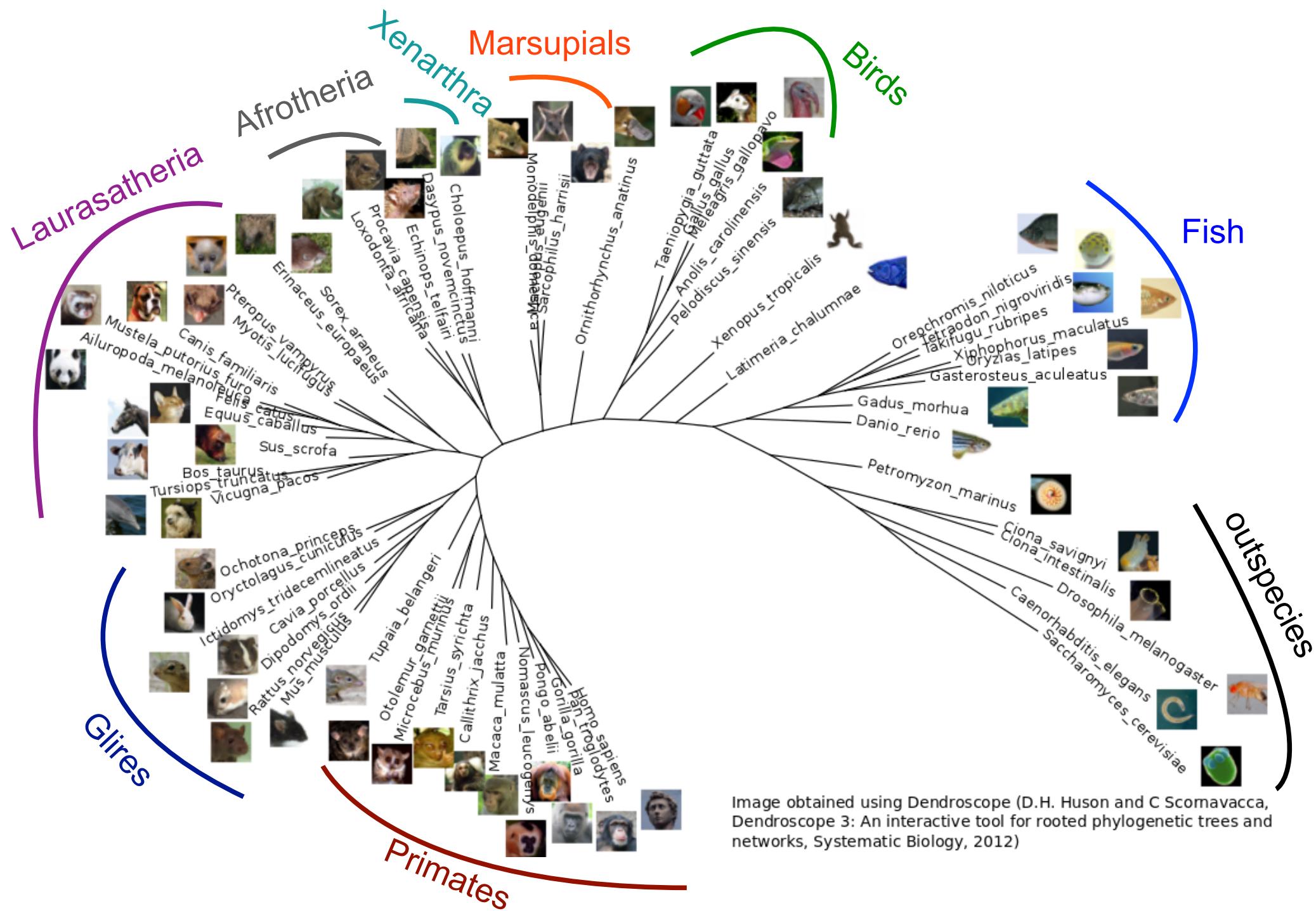
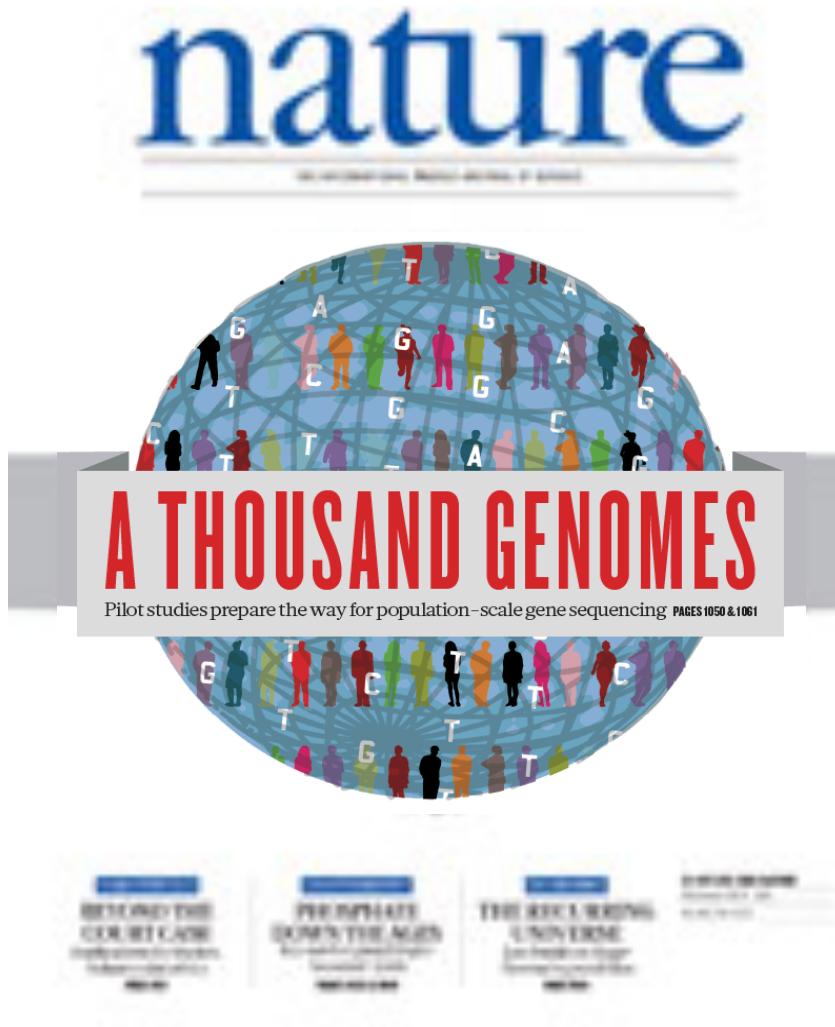


Image obtained using Dendroscope (D.H. Huson and C Scornavacca, Dendroscope 3: An interactive tool for rooted phylogenetic trees and networks, Systematic Biology, 2012)

1000 genomes pilot project (Oct 2010)



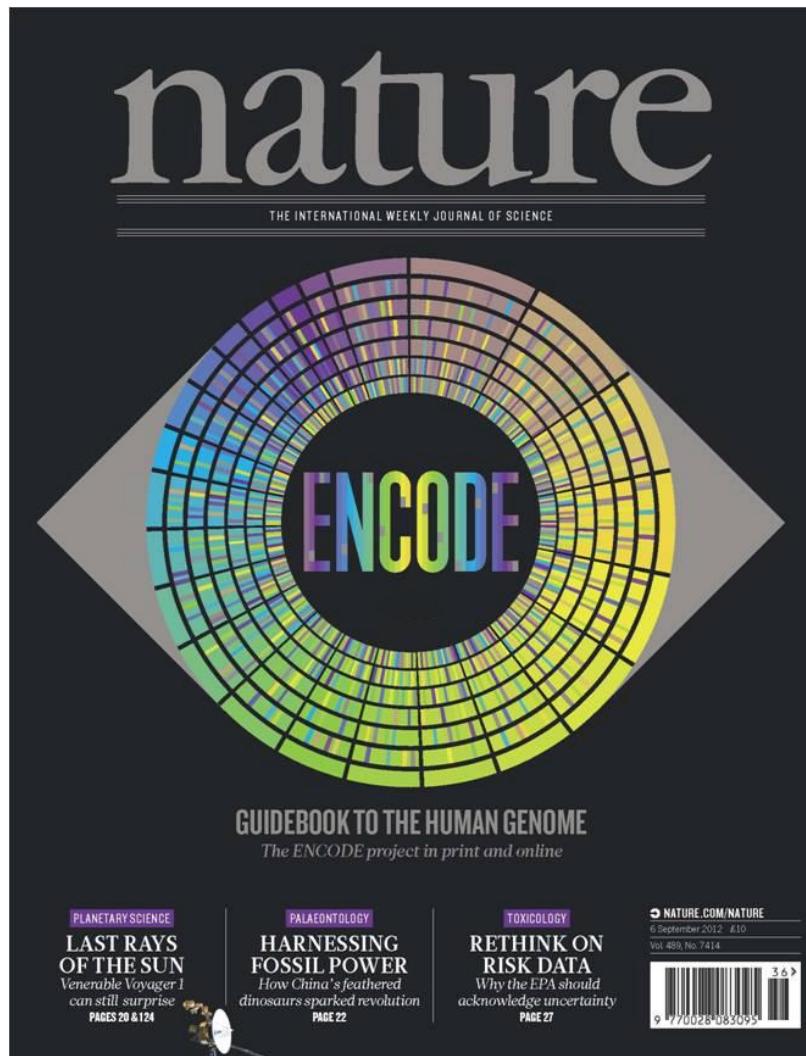
3 pilots:

- 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



- 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

Nature 489, 57–74 (06 Sept 2012)

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...

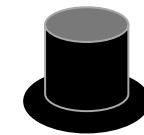
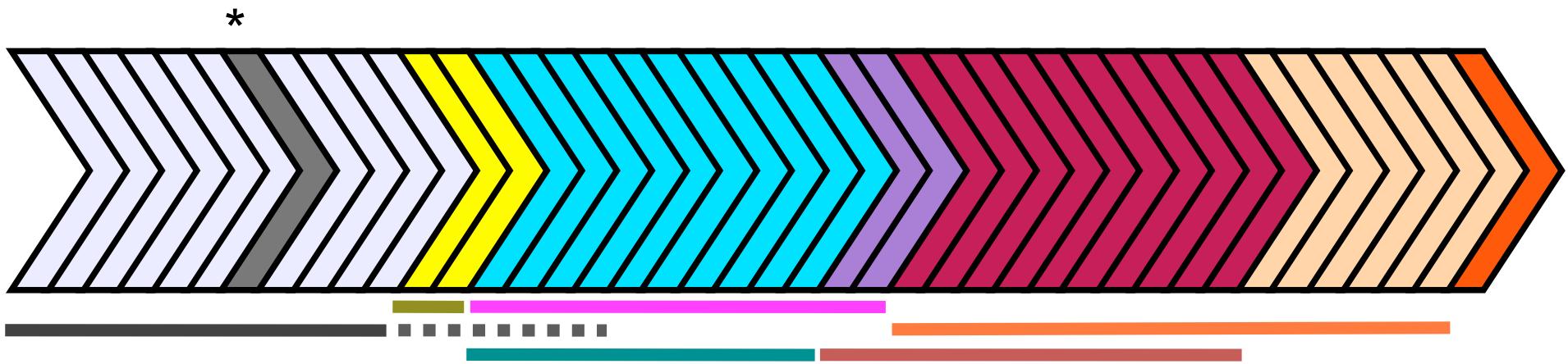
Variation, Comparative & Functional Genomics

- 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



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

D48-D55 *Nucleic Acids Research*, 2013, Vol. 41, Database issue
doi:10.1093/nar/gks1236

Published online 30 November 2012

Ensembl 2013

Paul Flicek^{1,2,*}, Ikhlaq Ahmed¹, M. Ridwan Amode², Daniel Barrell², Kathryn Beal¹, Simon Brent², Denise Carvalho-Silva¹, Peter Clapham², Guy Coates², Susan Fairley², Stephen Fitzgerald¹, Laurent Gil¹, Carlos García-Girón², Leo Gordon¹, Thibaut Hourlier², Sarah Hunt¹, Thomas Juettemann¹, Andreas K. Kähäri², Stephen Keenan¹, Monika Komorowska¹, Eugene Kulesha¹, Ian Longden¹, Thomas Maurel¹, William M. McLaren¹, Matthieu Muffato¹, Rishi Nag², Bert Overduin¹, Miguel Pignatelli¹, Bethan Pritchard², Emily Pritchard¹, Harpreet Singh Riat², Graham R. S. Ritchie¹, Magali Ruffier¹, Michael Schuster¹, Daniel Sheppard², Daniel Sobral¹, Kieron Taylor¹, Anja Thormann¹, Stephen Trevanion², Simon White², Steven P. Wilder¹, Bronwen L. Aken², Ewan Birney¹, Fiona Cunningham¹, Ian Dunham¹, Jennifer Harrow², Javier Herrero¹, Tim J. P. Hubbard², Nathan Johnson¹, Rhoda Kinsella¹, Anne Parker², Giulietta Spudich¹, Andy Yates¹, Amonida Zadissa² and Stephen M. J. Searle²

¹European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton Cambridge CB10 1SD, UK and

²Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SA, UK



European Commission
Framework Programme 7



Quantomics

From Sequence to Consequence :
Tools for the Exploitation of Livestock Genomes

