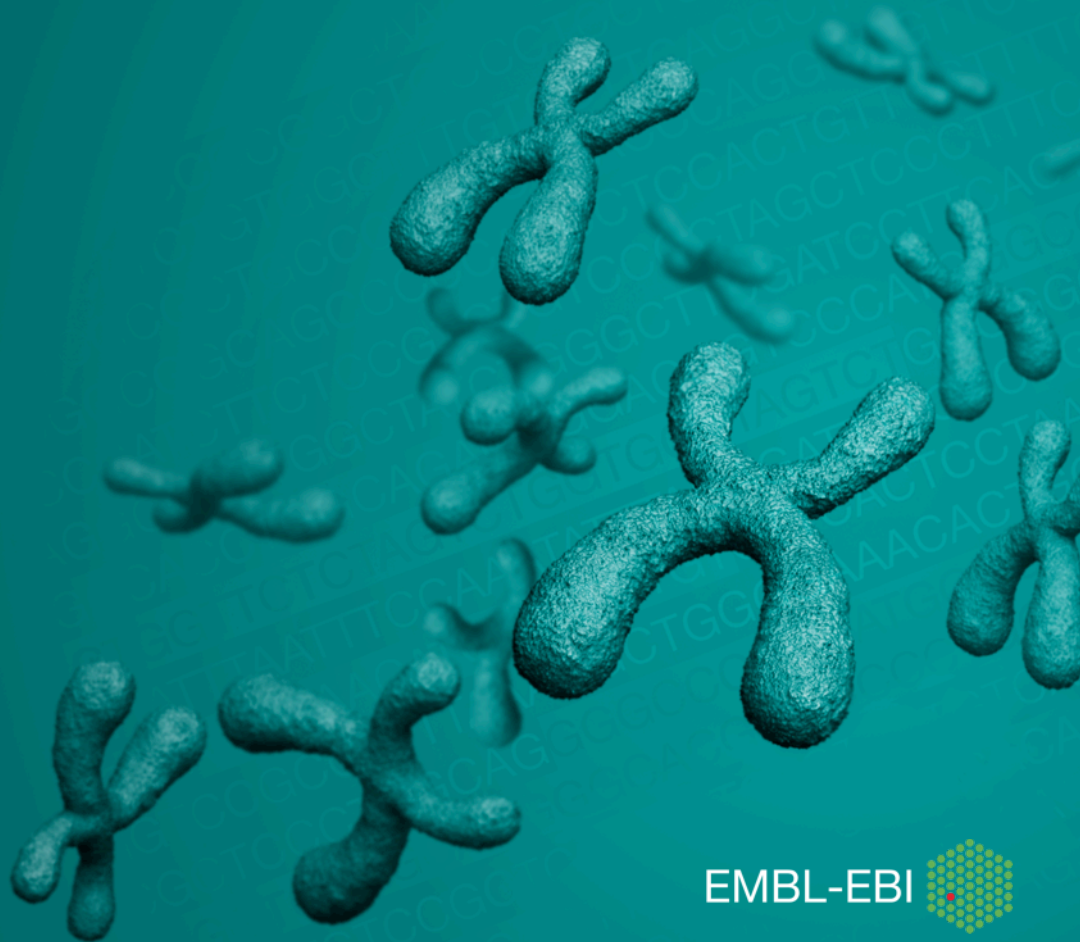# Genome Annotation

Ewan Birney (tweetable)
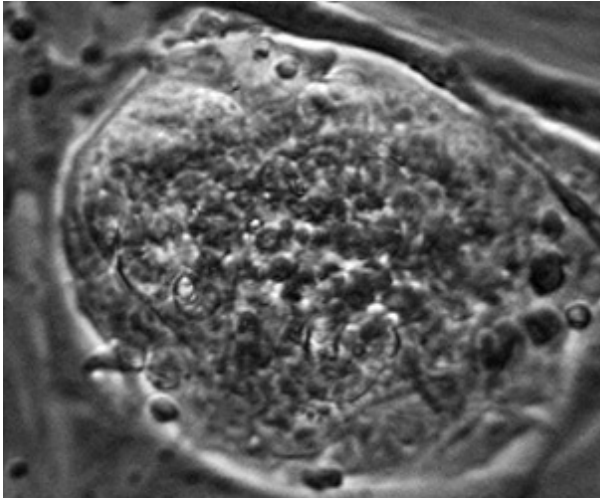
@ewanbirney

# What next?

```
GGGATTGAGAGTGATCACTCACGCTAACGTCTGCCCTGTTCCTGTATGGTGAGGCCGCAC
CACAAGCCACCACCGCCGCCGCCTTCTGCGCAACGCCAACCGCCCGCCAAAACGGATCCT
TCCCTGCGCCTGCGCAACCAATCTTGGGACCGGACCTTTTTTCTCCGCCCACTACGCATG
CGCAAAGCTAGGACAAACTCCCGCCAACACGCAGGCGCCGTAGGTTCACTGCCTACTCCT
GCCCGCCATTTCACGTGTTCTCAGAGGCAGGTGGAACTTCTTAATGCGCCTGCGCAAAAC
TCGCCATTTTACTACACGTGCGGTCAACAAGAGTTCATTGCAAAAAAATTGTTACCTCCT
AGCTGCTTGTCTAATACATAGTGTTAATCATGCTTTGCCAAGCGACTTGACTGTAATATT
TGCGCGTGGAAGATTAAAAAGATGTTAAACACCCAAGGTAGATTCAAATGTGAATGATTG
GTCGGTTGGCCAATCAGACTGGTTAACAATAACATTACTCGGGAACCAATGGACTCCAAG
GGGTGGAGACGGCGTAGAACGACCGAAGGAATGACGTTACACAGCAATGTGGCACCACAG
GCCAATAGCAGGGGGAAGCGATTTCAAGTATCCAATCAGAGCTGTTCTAGGGCGGAGTCT
ACCAATGCCGAAAGCGAGGAGGCGGGGTAAAAAGAGAGGGCGAAGGTAGGCTGGCAGAT
ACGTTCGTCAGCTTGCTCCTTTCTGCCCGTGGACGCCGCCGAAGAAGCATCGTTAAAGTC
TCTCTTCACCCTGCCGTCATGTCTAAGTCAGAGGTGAGTTAGGCGCGCTTTCCCACTTGA
ATTTTTTCCTCTCCCTTTCCTGAATCGGTAAGATGCTGCTGGGTTTCGTTCCTTGCACCA
GCCCATTCTACAGTTCCTTCGGTCGCTGCCACGGCCTACCCCTCCCAAAGTTCAAGTCGC
CATTTTGTCCTCTTGATCGCCATGAGGCCGCTCTCCGCCAACCATGTGTTATCATGCGGG
ACTCGTTACTCGTAGCAAAATTCTTAGGCACACAGGATCTTTGTCTTTTTTTAAACCTTG
CCTTGGTGAGCGAGTTTTCTAAAGAGCGATTAGTCCCATTGTGGAGATGCACCCCTACCG
CCCAAGCCTTTGTTGCGCGTGCGTCGGAAGGCGACTAGGGACGCATGCGCTTGCGATTTC
CTAGCACTCCCAACTCCAGCATACGGCCTCCCTTGATAGGCAGAAGCACGTGTCTTGTTG
CGACCTGAACGAACAATAAGTGCTAGGTACACAGTTGGTGTCTAGTTTTTCTTTTCCTCG
ATGGAAATTGTTTCGTGTTGTAGCCCATTTAACACTTCCCCCTCCCCCCACTCTAGTCTC
```
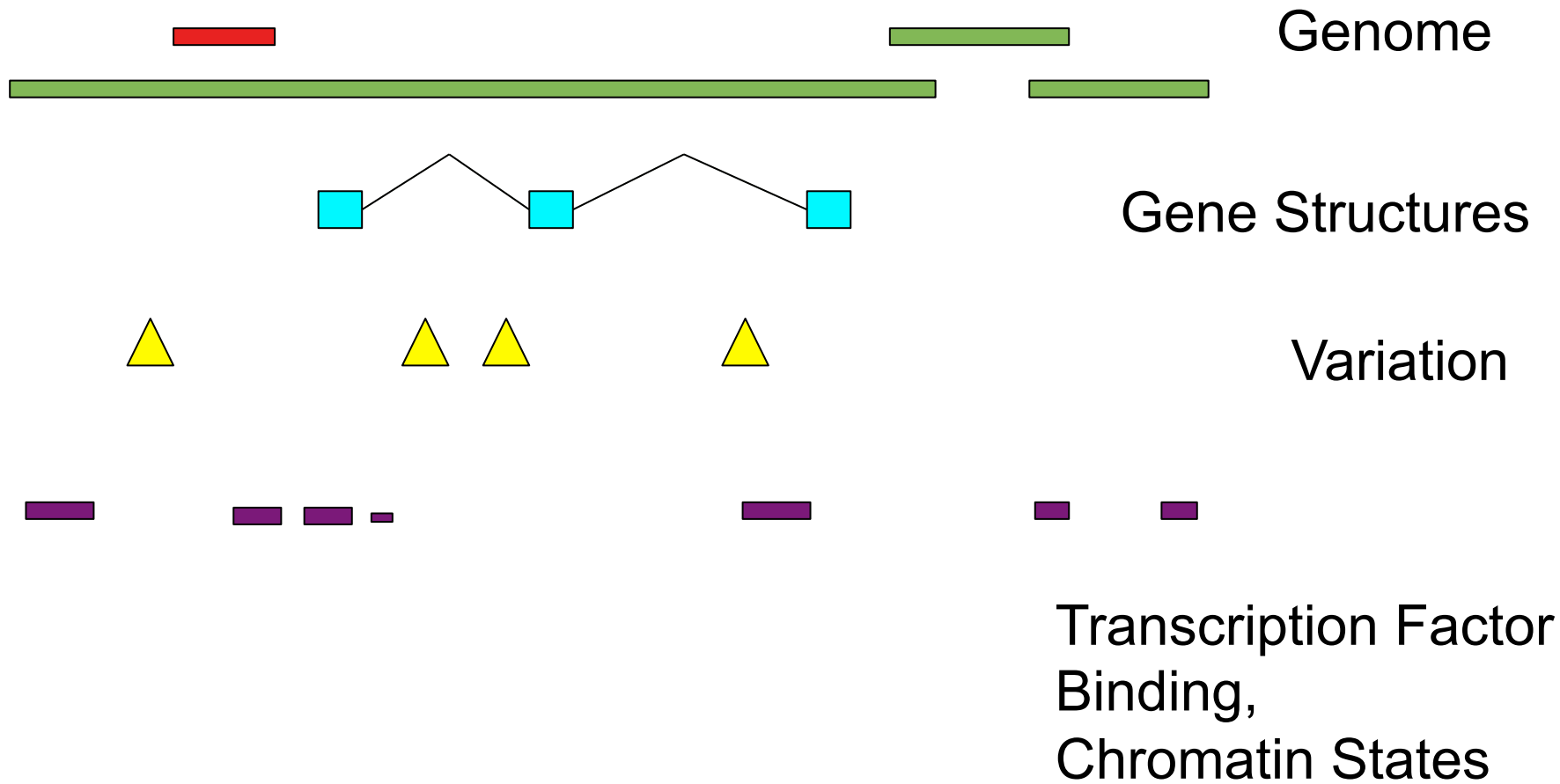
# The ultimate biological index



All biological molecules are encoded in some manner in the genome

…as is cellular response

…as is development

# Genome Annotation

Genome

Gene Structures

Variation

Transcription Factor Binding, Chromatin States

# Generation, Integration, Annotation

Generate                    Integrate                    Annotate


High quality data                                        Provide initial labeling
Good data stds                                           Understand (a bit more…)
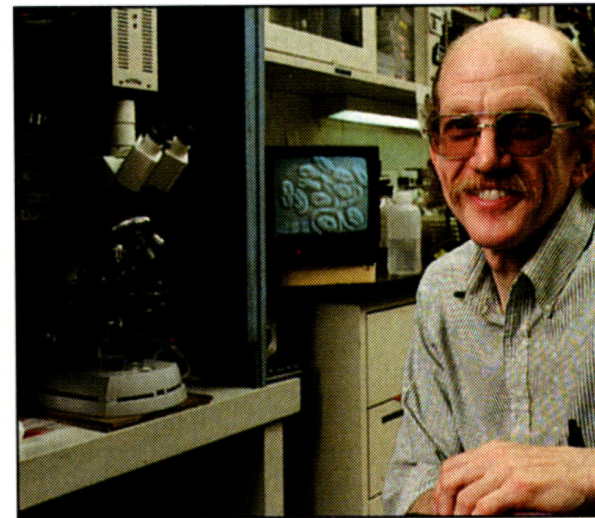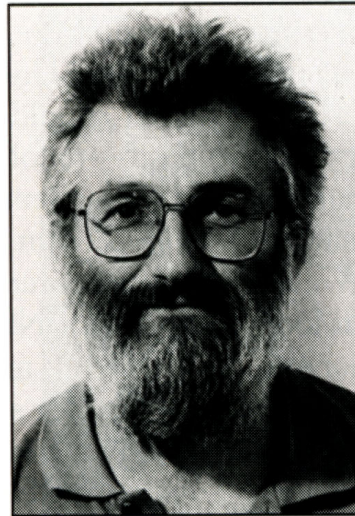Public availability(!)


Map to genome sequence
Make non redundant
Call

EMBL-EBI

# A Strategy for Sequencing the Genome 5 Years Early

In meetings over the past 6 weeks, two respected gene sequencers have been delivering a startling message: The chief goal of the Human Genome Project—obtaining a complete sequence of the 3 billion bases in human DNA—can be achieved as early as 2001, 5 years ahead of schedule. What is more, they say, it can be done without any fancy new technology. The two optimists—John Sulston, director of the Sanger Center in Cambridge, United Kingdom, and Robert Waterston, director of the Genome Sequencing Center at Washington University in St. Louis—have sketched a plan that they think could deliver the Holy Grail of genomics for $300 million to $400 million over 5 years. The U.S.
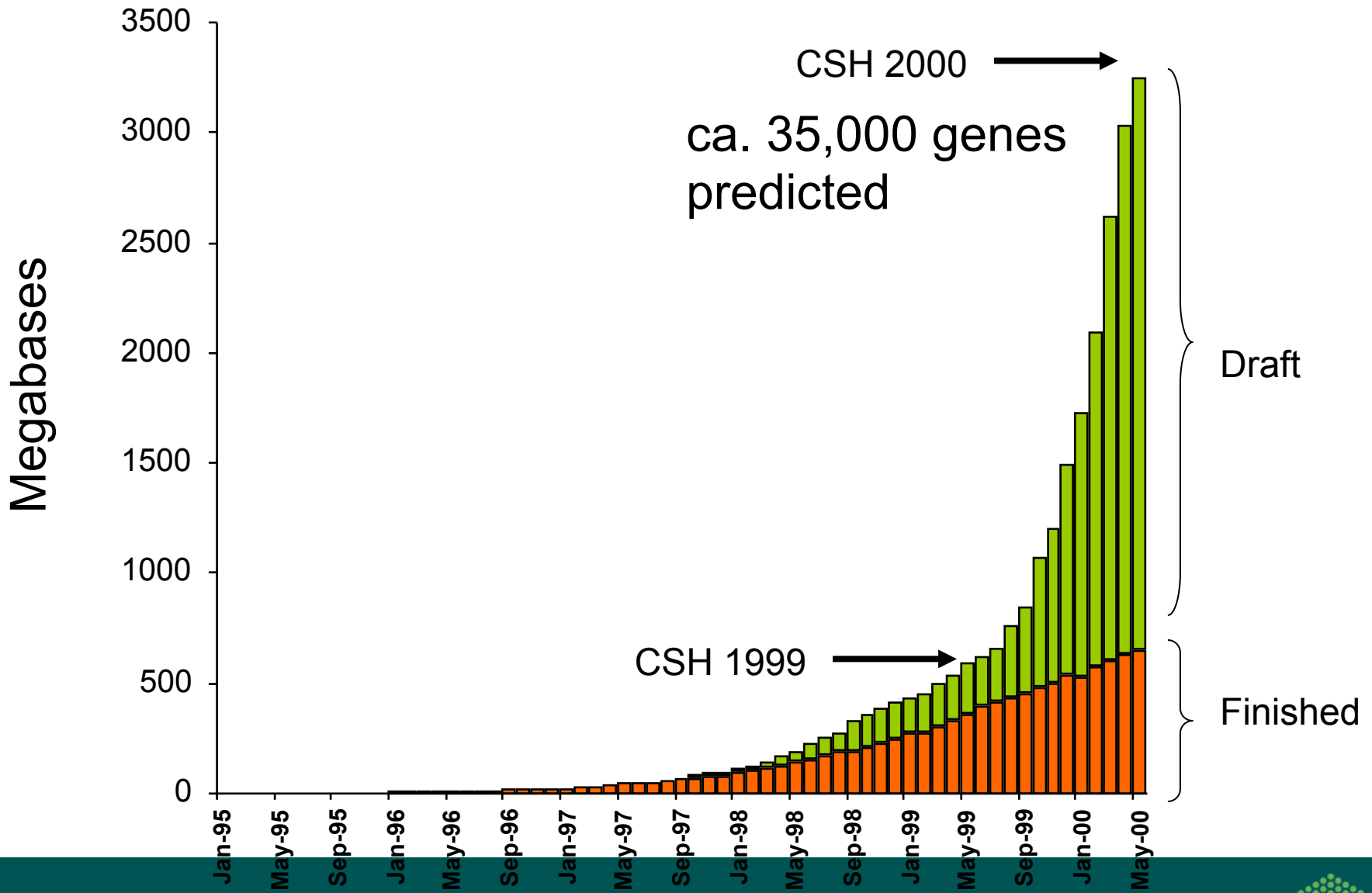
tion would lead to a breakthrough that would sharply cut the cost of obtaining sequence data. But that hasn't happened.

Now, Waterston and Sulston are arguing that instead of waiting for the ideal technology, it is time to move pragmatically into the final phase of the program—sequencing the
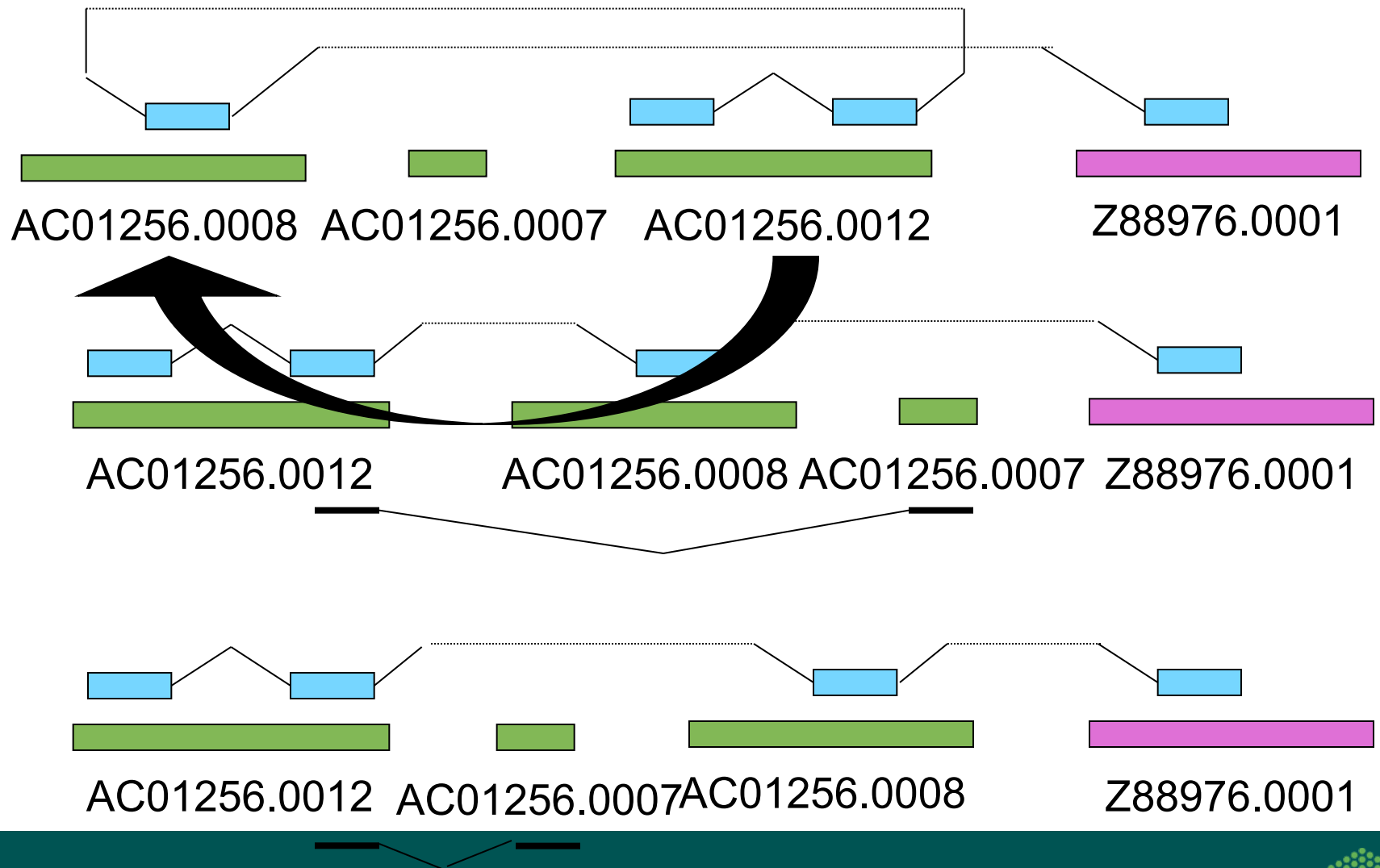
**Full speed.** John Sulston *(left)* and Robert Waterston have been floating to shift into high gear on sequencing, using current technology.

EMBL-EBI

# Human Genome Sequence in INSDC

# Ensembl (circa 2001!)



AC01256.0008  AC01256.0007  AC01256.0012  Z88976.0001

AC01256.0012  AC01256.0008 AC01256.0007  Z88976.0001

AC01256.0012  AC01256.0007 AC01256.0008  Z88976.0001

EMBL-EBI
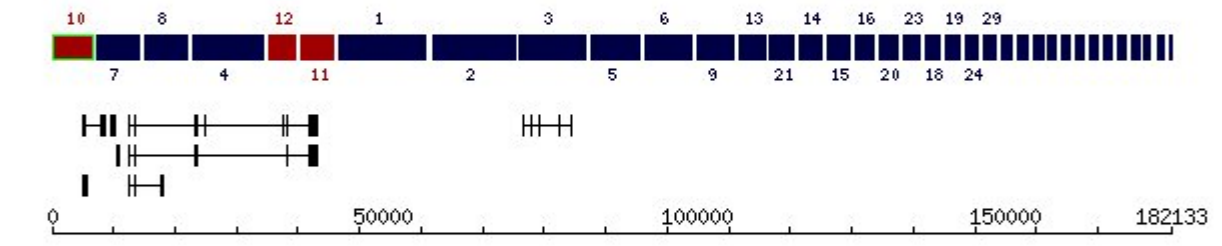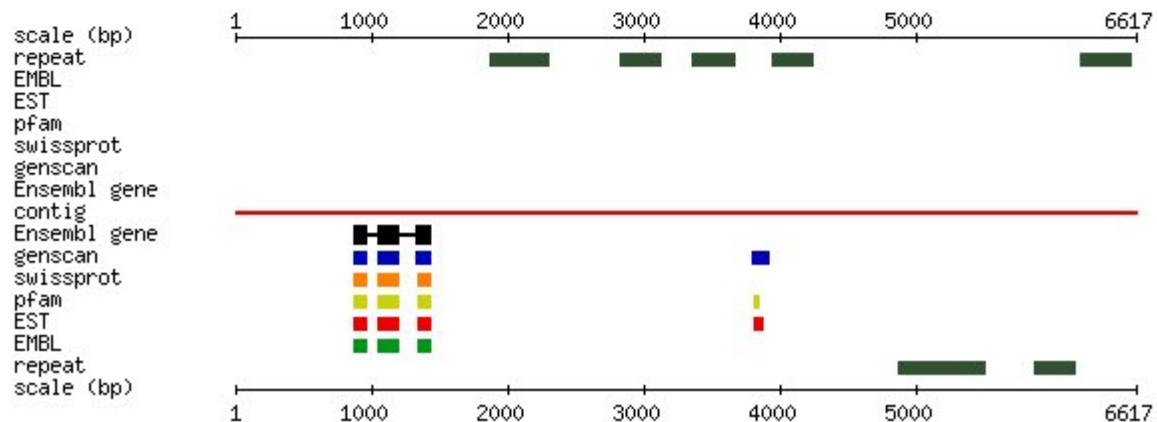
# Web Site Contig view



Click on a contig in the map below to view more detail or a transcript to see exon data and supporting evidence

Detailed Contig View : AP000869.00010

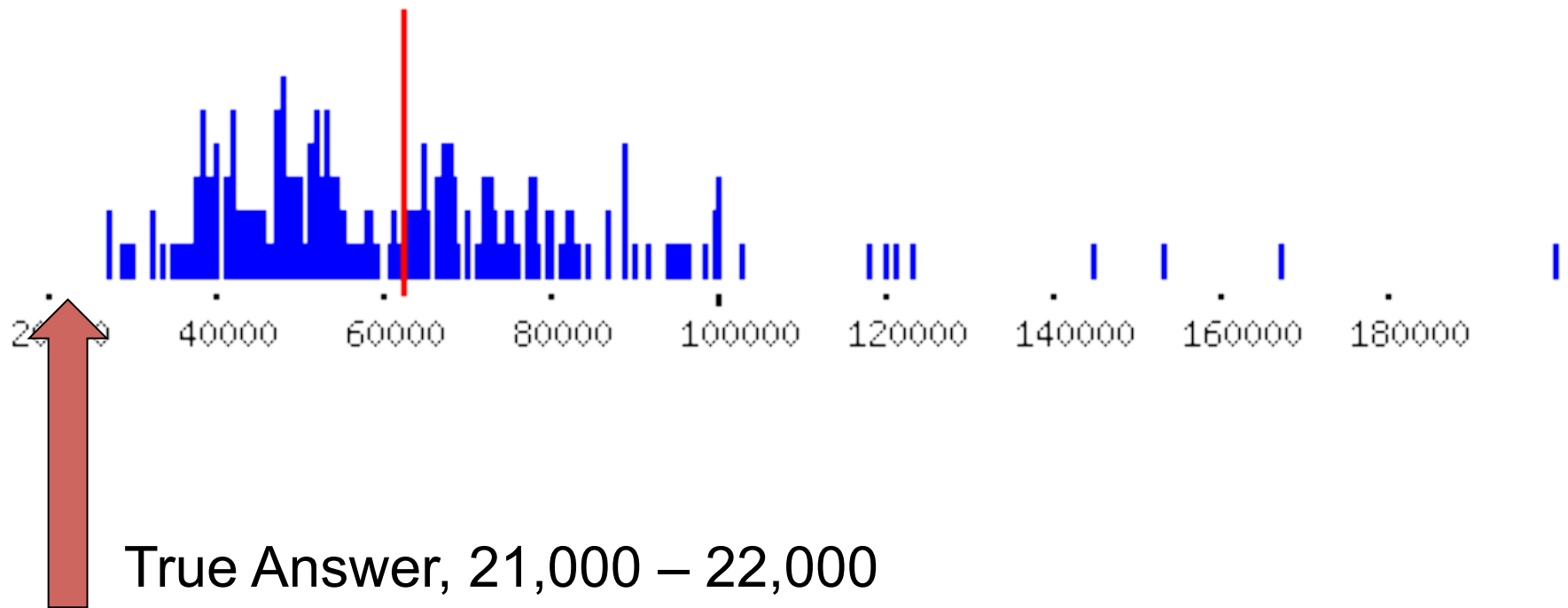This is a detailed view of the contig selected in the map above (highlighted with a green border)

Click a gene in the contig display below to view detailed gene information

# Estimated number of genes

- Chromosome 22
  - 549 +  (~100) genes
  - 1.1 % of genome
  - 50,000 genes (59,000)
- Chromosome 21
  - 225 genes in Chr21
  - 1.1 % of genome
  - 20,500 genes
- Ensembl
  - 38,000 genes

# Distribution of bets for the number of protein coding genes (2001)



True Answer, 21,000 – 22,000

EMBL-EBI

# Variation

Generate          Integrate          Annotate

SNP Consortium                      Frequency, LD
Perlgen                             VEP, VAAST
HapMap
1,000 Genomes

          dbSNP RefSNP
          1000 Genomes Pipeliens

# TFs and Chromatin

Generate                    Integrate                    Annotate

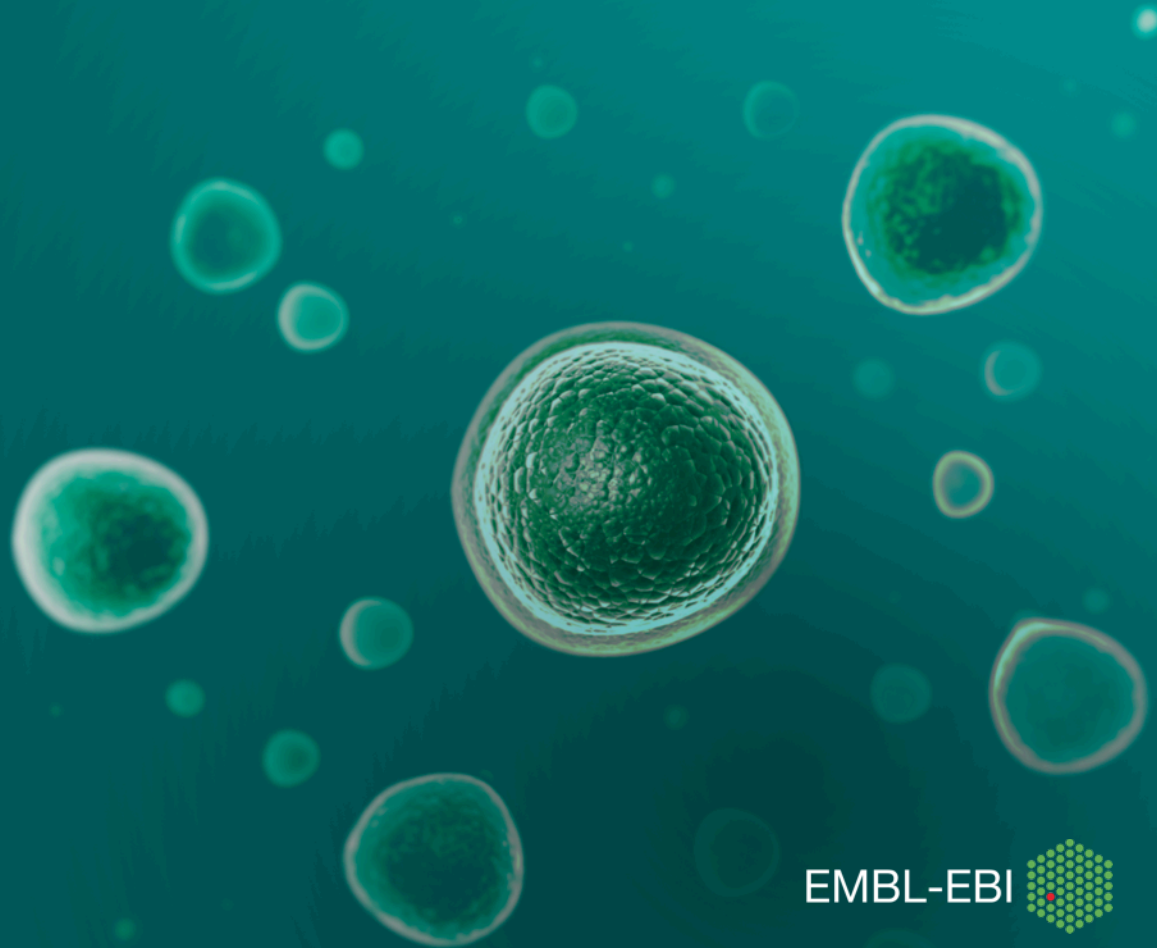ENCODE                                                   ? Chromatin State
Epigenome Roadmap                                        4C/5C Experiments?
IHEC
Individual Studies

ENCODE/Epigenome Pipelines
Ensembl Regulatory Build
RegulomeDB

EMBL-EBI

# ENCODE

EMBL-EBI

# ENCODE Dimensions



182 Cell Lines/ Tissues

Expression Array
RNA
Open chromatin
Histone Mods
TFs
Methylation

Cells

GM12878
H1-hESC
K562
HeLa-S3
HepG2
HUVEC
chr8

Genome

3,010 Experiments
5 TeraBases
1716x of the Human Genome

Methods/Factors

GM12878
K562
H1-hESC
HeLa-S3
HepG2
Huvec

Histone Mods
Pol2/3
Transcription Factors
Control

164 Assays (114 different Chip)

EMBL-EBI

# A consortium effort…

11 Main, multi Site groups
~50 Laboratories in total

10 additional groups

30 "lead" PIs

~410 Authors on the main Paper

6 "high profile" papers
~25-30 companion papers

Steve Landt @Stanford

Alexias Safi @Duke
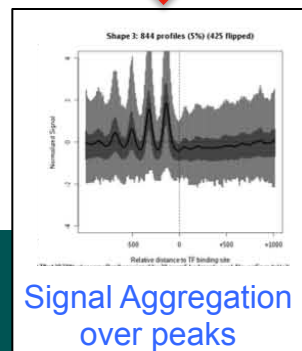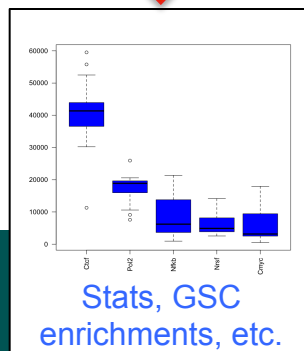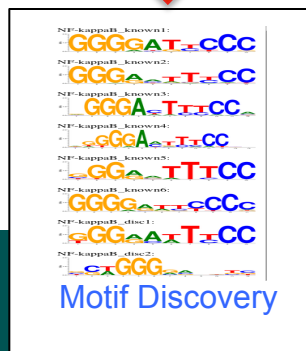
Chuck Epstein @Broad

Flo Pauli @HA

Jen Harrow @Sanger

Kate Rosenbloom@UCSC

Raj Kaul@UW

Carrie Davis@CSHL

EMBL-LBI

# ENCODE Uniform Analysis Pipeline

*Anshul Kundaje, Qunhua Li, Michael Hoffman, Jason Ernst, Joel Rozowsky, Pouya Kheradpour*



Mapped reads from production (Bam)

Uniform Peak Calling Pipeline (SPP, PeakSeq)

Signal Generation
(read extension and mappability correction)

Good reproducibility   Poor reproducibility

IDR Processing, QC and Blacklist Filtering

Segmentation

ChromHMM/Segway

BroadHistoneK562ControlStdAln_1Reps

Self Organising Maps

Motif Discovery

Stats, GSC enrichments, etc.

Signal Aggregation over peaks

EMBL-EBI

# Irreproducible Discovery Rate (IDR)

*Ben Brown, Qunhau Li, Peter Bickel*

If one re-ran the experiment, what is the probability one would observe the same element at this rank or better

Uses ranked element lists from two replicates, and makes the assumption that there is noise at the bottom of the rank



Chip-seq                     Dnase-seq                     RNA-seq

EMBL-EBI

# Raw genome coverage of elements

| Element Type | Coverage | Cumulative Coverage |
| --- | --- | --- |
| Exons | 3% | 3% |
| Chip-seq bound motifs | 4.5% | 5% |
| DNaseI Footprints | 5.7% | 9% |
| Chip-seq bound regions | 8.1% | 12% |
| DNaseI HS regions | 15.2% | 19.4% |
| Histone Modifications (*) | 44% | 49% |
| RNA | 62% | 80% |
| (* excluding broad marks) | | |

*(Union over all experiments and cell types)*

Region

Bound Motif/
Footprint

# Saturation

*Steve Wilder*



Most aggressive fit for saturation suggests a maximum of 50% of elements discovered

Likely to be lower due to inaccessible cell types etc

EMBL-EBI

# Evenly spaced over the genome

99% of the genome is within 1.7 KB of a biochemical event

95% of the genome is within 8 KB of a bound motif or footprint

EMBL-EBI

# Many other stories



**Splicing/Histone interaction (Roderic Guigo)**

**RNA landscape**
**Tom Gingeras**

**TF Co association,**
**Mike Snyder+Mark Gerstein**

**DNAseI footprints – John Stam.**
**DNA Methylation – Rick Myers**

# Discovering functional genome segments

*Michael Hoffman, Jason Ernst, Bill Noble, Manolis Kellis*



~7 Major flavours of genome
25 "elaborations"
1,000s of details

Well understood:
TSS, Gene Start,
Gene Bodies

Reassuringly Interesting
"Enhancers" (2 states)
Insulators

Definitely There, Unexpected
Specific Gene End

Sub-classification of Repeats

EMBL-EBI

# Fish Transgenics

- 7 strong positives, 6 negatives from Segmentation, 3 Blood

- 2 strong positives, 4 weak, 5 negatives from Naïve picks (Fisher's Exact 0.0393)

EMBL-EBI

# Many other stories



K562 Whole-genome

**TF Co association,**
**Mike Snyder+Mark Gerstein**



H3k9ac: upstream exon — 461 (+) in Gm12878–K562 — 180 (–) in Gm12878–K562 — p(bonferroni,g)=1.39

H3k9ac: AS exon — p(bonferroni,g)=0.0133

H3k9ac: downstream exon — p(bonferroni,g)=4.13

**Splicing/Histone interaction (Roderic Guigo)**

**RNA landscape**
**Tom Gingeras**



Gencode Annotation Features

**DNAseI footprints – John Stam.**
**DNA Methylation – Rick Myers**

EMBL-EBI

# Functional SNPs

*Belinda Giardine, Marc Shaub, Ross Hardison, Mike Snyder, John Stam.*

Genome Wide Association Studies (GWAS) Results

Linkage Disequilibrium

ENCODE Functional Region

Reported SNP

Statistically associated with the phenotype

fSNP

✔ Associated with the phenotype
✔ In a functional region

# Functional SNP - Direct Hit

Genome Wide Association
Studies (GWAS) Results

ENCODE Functional
Region

fSNP Direct Hit

✔ Association reported in a GWAS
✔ In a functional region

EMBL-EBI

# Direct hits

*Ross Hardison, Belinda Giardine, Marc Schaub*



1.3/1.2 enrichment vs matched null.

When you extend to SNPs in high LD, GWAS SNPs overlap by 80%

# Association of TF or Cell Type with Disease

*Ross Hardison, Belinda Giardine*



EMBL-EBI

# Zoom in…

| Phenotype | SNP-Pheno associations | overlap any TF occupancy | Gm12878Mef2a | Gm12878Pol2 | Gm12878Ebf | Gm12878Pol24h8 | Gm12878NfkbIggrab | Gm12878Irf4 | Gm12878Pax5c20 | Gm12878Pu1 | Gm12878Batf | HuvecGata2Ucd | Gm12878Egr1V0416101 | Helas3CebpbIggrab | Hepg2Ctcf | HUVEC.all | hTH1.all | hTH2.DS7842 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| TOTAL | 4860 | 600 | 47 | 69 | 78 | 57 | 35 | 35 | 54 | 47 | 45 | 29 | 28 | 69 | 54 | 85 | 192 | 57 |
| Multiple_sclerosis | 71 | 15 | 4 | 3 | 4 | 3 | 2 | 3 | 4 | 3 | 2 | 2 | 1 | 1 | 0 | 1 | 5 | 4 |
| Systemic_lupus_erythematosus | 62 | 10 | 4 | 6 | 4 | 6 | 4 | 3 | 1 | 1 | 4 | 1 | 2 | 2 | 1 | 2 | 4 | 2 |
| Height | 204 | 34 | 3 | 3 | 7 | 3 | 2 | 5 | 3 | 1 | 5 | 0 | 6 | 7 | 6 | 6 | 9 | 3 |
| Rheumatoid_arthritis | 57 | 11 | 3 | 2 | 4 | 2 | 4 | 0 | 4 | 4 | 2 | 0 | 1 | 1 | 0 | 2 | 11 | 3 |
| Chronic_lymphocytic_leukemia | 17 | 8 | 1 | 5 | 1 | 4 | 1 | 1 | 2 | 3 | 1 | 0 | 0 | 0 | 0 | 0 | 2 | 0 |
| Celiac_disease | 54 | 11 | 1 | 3 | 4 | 3 | 2 | 1 | 1 | 2 | 1 | 0 | 0 | 0 | 0 | 0 | 2 | 1 |
| Ulcerative_colitis | 85 | 11 | 3 | 3 | 2 | 3 | 3 | 1 | 3 | 2 | 2 | 2 | 2 | 0 | 1 | 2 | 7 | 2 |
| Crohn's_disease | 105 | 20 | 2 | 2 | 2 | 2 | 1 | 1 | 2 | 2 | 1 | 5 | 1 | 2 | 2 | 6 | 9 | 5 |

EMBL-EBI

# Example loci

*Ross Hardison, Belinda Giardine*

,819,500 (1,545,000 bp)

40500000

*C9*

*DAB2*

*BC026261*

*PTGER4*

*TTC33*

*OSRF*

*PRKAA1*

GWAS Catalog

chr5:40,390,001-40,440,000 (50,000 bp)

- ● Crohn's disease
- ● ulcerative colitis
- ● multiple sclerosis

rs4613763

rs17234657

rs11742570

rs1992660

rs6451493

rs6896969

rs1373692

rs9292777

TFs
- HUVEC GATA2
- HUVEC cFOS
- HUVEC Input

DNase I
- HUVEC
- Jurkat
- Th1
- Th2

EMBL-EBI

# Organisation+Access for the data

## www.encodeproject.org (UCSC)

"Factorbook"



### UCSC Genome Browser

### Ensembl





Raw Data in GEO/SRA

EMBL-EBI

# Ensembl Regulation



- Integration over ENCODE, NIH Epigenome Roadmap, IHEC, some individual lab

- "Experiment", not per-replicate level view, only QC passed data

- Progressive integration into "one track" on the genome

- Integration into VEP - variant effect predictor (with SIFT and POLYPHEN as well!)

- Future:
  - Linking of genes to regulatory element, Linking of tfs to phenotypes

- Ensembl workshops (we come to you)

- Ensembl course (you come to us)

- helpdesk@ensembl.org

EMBL-EBI

# And… just for fun…

# Over a beer…

Ha! At some point all the data we
Store is going to be DNA…

Of course, the cost effective way
To store this would be as DNA…

**Figure 2 | Digital information encoded in DNA.** Digital information (**A**, in blue), here binary digits holding the ASCII codes for part of Shakespeare's sonnet 18, was converted to base-3 (**B**, red) using a Huffman code. This in turn was converted *in silico* to our DNA code (**C**, green), with no homopolymers, which formed the basis for a large number of overlapping DNA segments each containing 100 bases of encoded information (**D**, green or, with alternate segments reverse complemented for added data security, violet) and with orientation and indexing DNA codes added (yellow, as described in the text). These strings were synthesised, sequenced and decoded. **E**, A digital photograph of the EMBL-European Bioinformatics Institute (JPEG 2000 format) and **F**, an extract of the Watson and Crick (1953) paper[10] (PDF format) that were among the files encoded in DNA and successfully recovered in this study.

EMBL-EBI

# Cost effective?

Dave Simonds

EMBL-EBI

# ENCODE Authors

**Ian Dunham, Anshul Kundaje**

**Shelley F. Aldred Patrick J. Collins, Carrie A. Davis, Francis Doyle, Charles B. Epstein, Seth Frietze, Jennifer Harrow, Vishwanath R. Iyer, Rajinder Kaul, Jainab Khatun, Bryan R. Lajoie, Stephen G. Landt, Bum-Kyu Lee, Florencia Pauli, Kate R. Rosenbloom, Peter Sabo, Alexias Safi, Amartya Sanyal, Noam Shoresh, Jeremy M. Simon, Lingyun Song, Nathan D. Trinklein**

**Funded by NIH-NHGRI**

*** Overall Coordination ***
Ian Dunham (1), Anshul Kundaje (2).

*** Data Production Leads ***
Shelley F. Aldred (3), Patrick J. Collins (3), Carrie A. Davis (4), Francis Doyle (5), Charles B. Epstein (6), Seth Frietze (7), Jennifer Harrow (8), Vishwanath R. Iyer (9), Rajinder Kaul (10), Jainab Khatun (11), Bryan R. Lajoie (12), Stephen G. Landt (13), B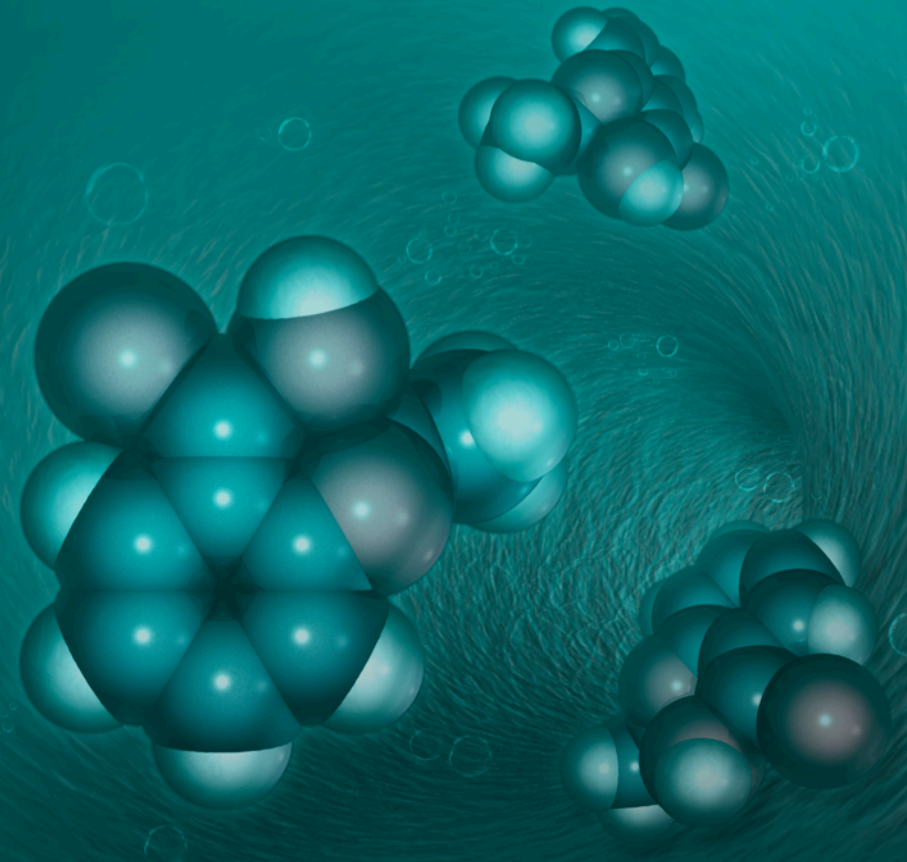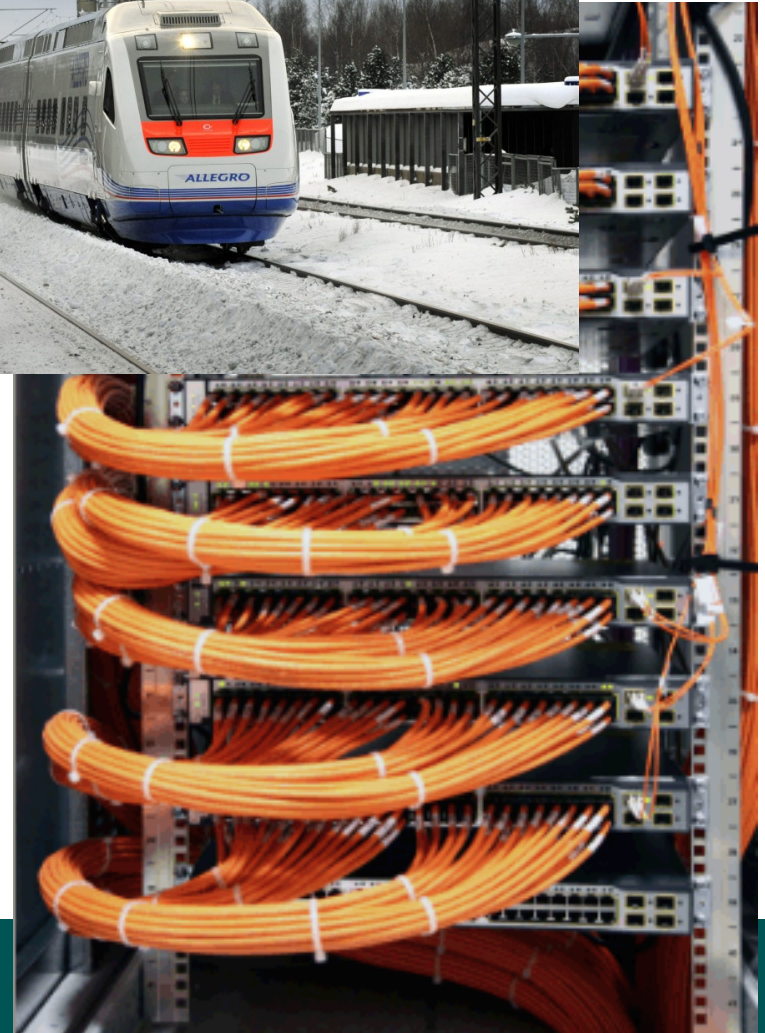um-Kyu Lee (9), Florencia Pauli (14), Kate R. Rosenbloom (15), Peter Sabo (16), Alexias Safi (17), Amartya Sanyal (12), Noam Shoresh (6), Jeremy M. Simon (18), Lingyun Song (17), Nathan D. Trinklein (3).

*** Lead Analysts ***
Robert C. Altshuler (19), E... 
Ian Dunham (1), Jason Ern...
Ross C. Hardison (26), Rol...
Kellis (19), Jainab Khatun...
Xinying Lin (23), Angelika...
Rozowsky (21), Felix Schle...
Steven P. Wilder (1), Weish...

*** Writing Group ***
Bradley E. Bernstein (33), ...

*** NHGRI Project Manager...
Leslie B. Adams (36), Laur...
Kelly (36), Rebecca F. Low...

*** Principal Investigators ...
Bradley E. Bernstein (33), ...
Farnham (7), Mark Gerstei...
Guig<F3> (40), Ross C. Ha...
(18), Elliott H. Margulies (2...
Tenenbaum (5), Zhiping W...

*** Broad Institute Group **...
Bradley E. Bernstein (33), ...
Mikkelsen (6), Shawn Gillespie (43), Alon Goren (33), Oren Ram (33), Xiaolan Zhang (6), Li Wang (6), ...
Michael J. Coyne (6), Timothy Durham (6), Manching Ku (33), Thanh Truong (6), Lucas D. Ward (19), Robert ...
(19), Matthew Eaton (19), Manolis Kellis (19).

*** Boise State University Proteomics Group ***
Jainab Khatun (11), Yanbao Yu (44), John Wrobel (11), Brian A. Risk (11), Harsha Gunawardena (44), Heather ...
(44), Christopher W. Maier (44), Ling Xie (44), Xian Chen (44), Morgan C. Giddings (11).

*** Data Coordination ...
Katrina Learned (15), Venkat S. Malladi (15), Kate R. Rosenbloom (15), Cricket A. Sloan (15), Matthew C. Wong (15), Galt ...
P. Barber (15), Melissa S. Cline ...
Kent (15), Vanessa M. Kirkup (15), Laurence R. Meyer, Jerry ... Morgan Maddren ...
(15).

*** Sanger Institute, Washington University, Yale University, Center for Genomic Regulation, Barcelona, ...
University of Lausanne, CNIO Group ***
Bronwen Aken (8), Roger P. Alexander (21), Suganthi Balasubramanian (21), Daniel Barrell (8), Gemma Barson ...
Andrew Berry (8), Nitin Bhardwaj (21), Alexandra Bignell (8), Veronika Boychenko (8), Michael Brent (45), Giovanni ...
Bussotti (22), Chao Cheng (21), Jacqueline Chrast (46), Claire Davidson (8), Thomas Derrien (47), Gloria ...
(8), Mark Diekhans (48), Jason Ernst (19), Iakes Ezkurdia (49), Julio Fernandez Banet (8), Adam Frankish ...
Gerstein (21), James Gilbert (8), Jose Manuel Gonzalez (8), Ed Griffiths (8), Roderic Guig<F3> (40), Lukas ...
Jennifer Harrow (8), Rachel Harte (48), David Haussler (50), C<E9>dric Howald (51), Timothy J. Hubbard ...
(8), Mike Kay (8), Manolis Kellis (19), Pouya Kheradpour (21), j Khurana (21), Felix Kokocinski (8), Jing ...
Michael F. Lin (19), Lucas Lochovsky (21), Jane Loveland (8), Zhi Lu (52), Deepa Manthravadi (8), Marco ...
Renqiang Min (21), Xinmeng (Jasmine) Mu (21), Jonathan Mudge (8), Gaurab Mukherjee (8), Cedric Notredame ...
Baikang Pei (21), Alexandre Reymond (21), Jose Manuel Rodriguez (49), Joel Rozowsky (21), Gary Saunders ...
Sboner (53), Stephen Searle (8), Cristina Sisu (21), Catherine Snow (8), Charlie Steward (8), Andrea Tanzer ...
Tapanari (8), Michael L. Tress (49), Alfonso Valencia (49), Marijke J. van Baren (55), Nathalie Walters (46) ...
Wilming (8), Koon-Kiu Yan (21), Kevin (Yuk-Lap) Yip (21), Amonida Zadissa (8), Zhengdong Zhang (56), Al...
(21), Alexej Abyzov (21).

*** Genome Institute of Singapore Group ***

*** HudsonAlpha ...
Devin M. Absher ...
Bowling (14), Ma...
(14), DeSalvo Gil...
Levy (14), Max W...
(30), Michael A. ...
Newberry (14), S...
Barbara Pusey (...
Vielmetter (42), E...

*** University of ...
Seth Bekiranov (12), Gaurav Jain (12), Bryan R. Lajoie (12), Amartya Sanyal (12).

*** University of Massachusetts Medical School Weng Group ***
Zhiping Weng (23), Troy W. Whitfield (23), Jie Wang (23), Patrick J. Collins (3), Shelley F. Aldred (3), ...
Trinklein (3), E. Christopher Partridge (14), Richard M. Myers (14).

Michael Snyder (35), Kevin P. White (41).

*** University of Heidelberg Group ***
Nathan D. ...
Thomas Auer (85), Lazaro Centanin (85), Michael Eichenlaub (85), Franziska Gruhl (85), Stephen Heermann (85), Daigo Inoue ...
Sinn

David Gonzalez (22), Assar Gordon (4), Harsha Gunawardena (44), C<E9>dric Howald ...
(47), Philipp Kapranov (68), Brandon King (23), Colin Kingswood (70), Guoliang Li ...
Luo (57), Eddie Park (30), Jonathan B. Preall (4), Kimberly Presaud (4), Paolo Ribeca (67), Brian A. ...
Robyr (72), Xiaoan Ruan (57), Michael Sammeth (67), Kuljeet Singh Sandhu (57), Lorain Schaeffer ...
See (4), Atif Shahab (57), Jorgen Skancke (22), Ana Maria Suzuki (66), Hazuki Takahashi (66), Hagen ...
Diane Trout (59), Nathalie Walters (46), Huaien Wang (4), John Wrobel (11), Yanbao Yu (44), Yoshihide Hayashizaki ...
(66), Jennifer Harrow (8), Mark Gerstein (21), Timothy J. Hubbard (8), Alexandre Reymond (51), Stylianos ...
Antonarakis (72), Gregory J. Hannon (4), Morgan C. Giddings (11), Yijun Ruan (57), Barbara Wold (42) ...
Carninci (66), Roderic Guig<F3> (40), Thomas R. Gingeras (39).

*** University of ...
Ha... Gaurav Jain (16), Anshul Kundaje (16), Rajinder Kaul (10) ...
Kris... Lee (16), Patrick Nease (16), Shane Joseph (16), Florence V. Nelson (16), Alex Reynolds (16) ...

*** SUNY ...
Anthony O. Shafer (16), George Stamatoyannopoulos (75), John A. Stamatoyannopoulos (16), Sean ...
Richard S. Sandstrom (16), Daniel L. Bates (16), Theresa K. Canfield (16), Amartya Sanyal (12) ...
Bob Thurman (16), Shinny Vong (16), Hao Wang (16), Molly A. Weaver (16).

*** Stanford-Yale, Harvard, University of Massachusetts Medical School, University of Southern California ...
Group ***
Alexej Abyzov (21), Nick Addleman (13), Roger P. Alexander (21), Raymond K. Auerbach (76), Suganthi ...
Balasubramanian (21), Keith Bettinger (13), Nitin Bhardwaj (21), Alan P. Boyle (13), Alina R. Cao (77), Philip Cayting ...
(13), Alexandra Charos (78), Chao Cheng (21), Yong Cheng (13), Catharine Eastman (13), Ghia Euskirchen (13) ...
Peggy Farnham (7), Joseph D. Fleming (79), Seth Frietze (7), Mark Gerstein (21), Fabian Grubert (13), Lukas ...
Karsten (13), Maya Kasowski (21), j Khurana (21), Phil Lacroute (13), Hugo Lam (13), Nathan Lamarre-Vincent ...
Stephen G. Landt (13), Jing (Jane) Leng (21), Jin Lian (82), Marianne Lindahl-Allen (79), Lucas Lochovsky ...
(21), Renqiang Min (21), Benoit Miotto (79), Hannah Monahan (79), Zarmik Moqtaderi (79), Xinmeng ...
(Jasmine) Mu (21), Henriette O'Geen (13), Zhengqing Ouyang (13), Dorrelyn Patacsil (13), Baikang Pei (21) ...
Debasish Raha (78), Lucia Ramirez (13), Brian Reed (78), Joel Rozowsky (21), Andrea Sboner (13), Minyi Shi (13) ...
Cristina Sisu (21), Teri Slifer (13), Michael Snyder (35), Kevin Struhl (79), Sherman M. Weissman (83), Heather Witt ...
(83), Linfeng Wu (13), Xiaoqin Xu (77), Koon-Kiu Yan (21), Xinqiong Yang (13), Kevin (Yuk-Lap) Yip (21), Zhengdong ...

William Stafford Noble (89), Jeffrey A. Bilmes (90), Orion J. Buske (16), Michael M. Hoffman (16), Daniel ...
Risk (11), Lucas Lochovsky (21), Peter V. Kharchenko (91), Peter J. Park ...
(91), Zhiping Weng (23), Sowmya Iyer (27), Xianjun Dong (23), Melissa Greven (23), Xinying Lin (23), Jie Wang (23), Hualin S. ...
Xi (32), Jiali Zhuang (23), Alexej Abyzov ...
(21), Raymond K. Alexander (21), Suganthi ...
Balasubramanian (21), Nitin Bhardwaj (21) ...
Chao Cheng (21), Lukas Habegger (21) ...
j Khurana (21), Jing ...
(Jane) Leng (21), Lucas Lochovsky (21) ...
Zhi Lu (52), Renqiang Min (21), Xinmeng ...
(Jasmine) Mu (21), Baikang Pei (21) ...
Andrea Sboner (53), Cristina Sisu (21) ...
Koon-Kiu Yan (21), Kevin (Yuk-Lap) Yip ...
(21), Mark Gerstein (21), Joel Rozowsky ...
(21), Zhengdong Zhang (56), Ewan Birney ...
(1).

EMBL-EBI

# Why we need a infrastructure

EMBL-EBI

# Infrastructures are critical…

# But we only notice them when they go wrong

# ELIXIR's mission

To build a sustainable European infrastructure for biological information, supporting life science research and its translation to:

society

bioindustries

environment
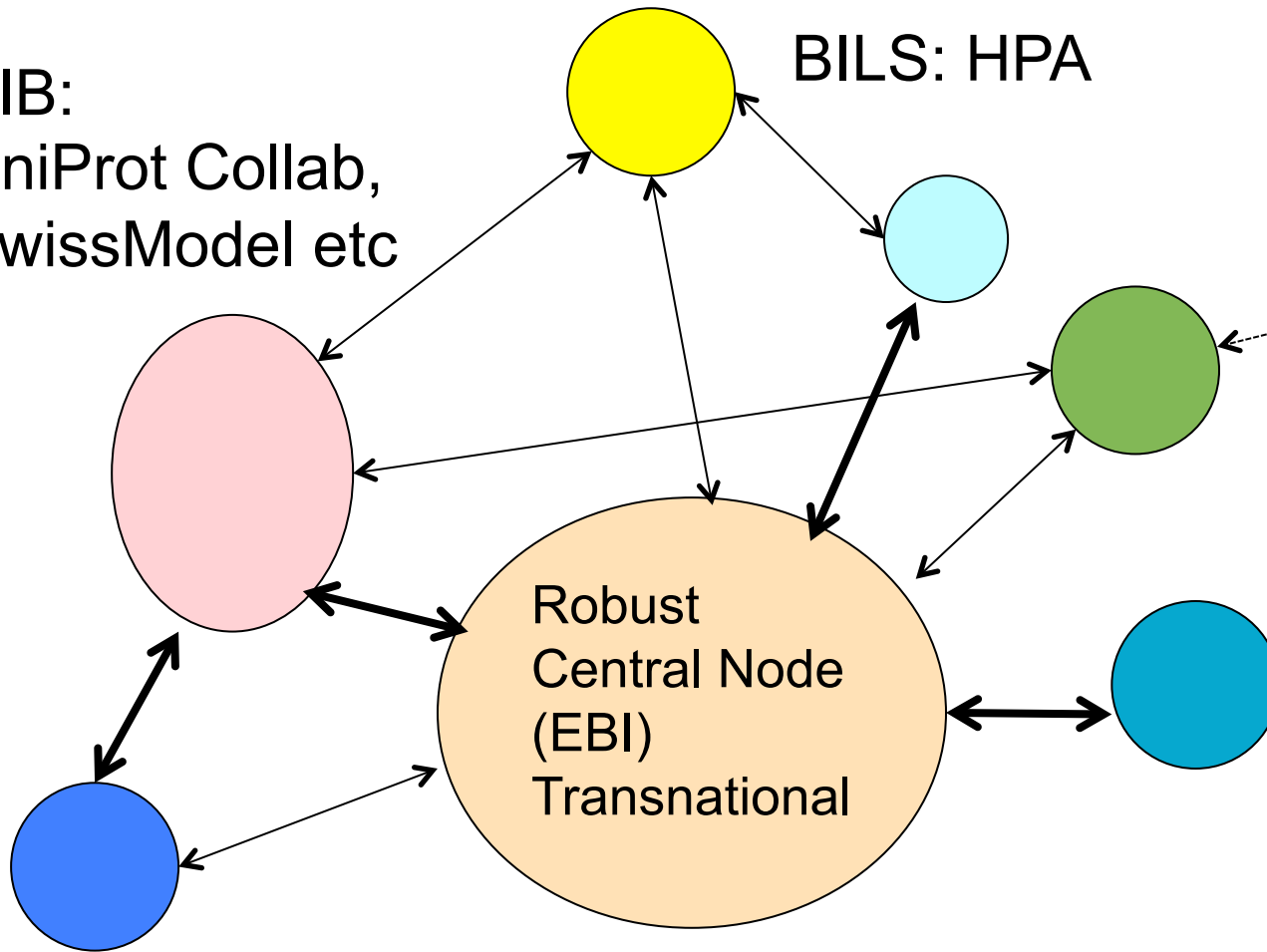
medicine

EMBL-EBI

# Other infrastructures needed for biology

- EuroBioImaging
  - Cellular and whole organism Imaging
- BioBanks (BBMRI)
  - We need numbers – European populations – in particular for rare diseases, but also for specific sub types of common disease
- Mouse models and phenotypes (Infrafrontier)
  - A baseline set of knockouts and phenotypes in our most tractable mammalian model
  - (it's hard to *prove* something in human)
- Robust molecular assays in a clinical setting (EATRIS)
  - The ability to reliably use state of the art molecular techniques in a clinical research setting

# Questions?

(you can follow me on twitter @ewanbirney)
I blog and update this on Google Plus publically

EMBL-EBI