

Lecture 15. Functional Genomics 4

Michael Schatz

March 30, 2017

JHU 600.649: Applied Comparative Genomics



Project Proposal

Due: Thursday March 30 @ 11:59pm

The screenshot shows a web browser window with three tabs open, all titled 'appliedgenomics/projectpropo'. The URL in the address bar is <https://github.com/schatzlab/appliedgenomics/blob/master/assignments/projects/projectproposal.md>. The browser interface includes a toolbar with various icons like Mail, JHUMail, Daily, and social media links. The GitHub header shows the user 'Michael' and navigation links for 'This repository', 'Search', 'Pull requests', 'Issues', and 'Gist'. Below the header, the repository name 'schatzlab / appliedgenomics' is displayed, along with statistics: 6 unwatched, 9 stars, and 0 forks. A 'Code' tab is selected. The main content area shows a single commit by 'mschatz' updating 'projectproposal.md'. The commit message is 'Update projectproposal.md' and it was made 4 minutes ago. The file details show 24 lines (16 sloc) and 1.46 KB. Below the file content, there is a section titled 'Project Proposal' containing assignment instructions and a list of proposal components.

Project Proposal

Assignment Date: March 16, 2017
Due Date: Thursday, March 30, 2017 @ 11:59pm

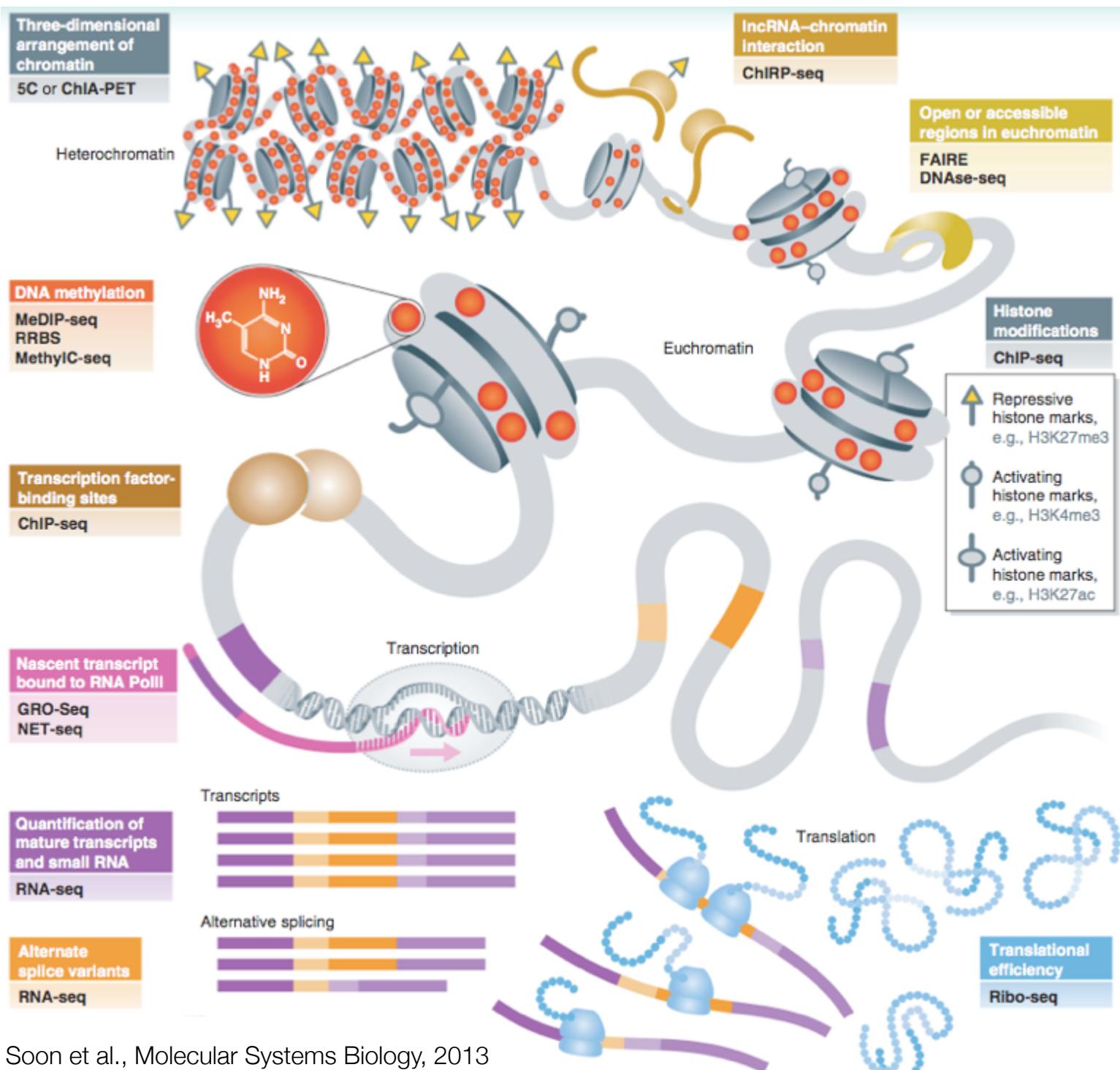
Review the [Project Ideas](#) page

Form a team for your class project (no more than 3 people to a team) and email a PDF of your project proposal (1/2 to 1 page) to "jhuappliedgenomics at gmail dot com" by 11:59pm on March 30.

The proposal should have the following components:

- Name of your team
- List of team members and email addresses
- Short title for your proposal
- 1 paragraph description of what you hope to do and how you will do it
- References to relevant papers
- References/URLs to datasets that you will be studying (Note you can also use simulated data)

After submitting your proposal, we will schedule a time to discuss your proposal, especially to ensure you have access to the data that you need. The sooner that you submit your proposal, the sooner we can schedule the meeting. No late days can be used for this project.



ARTICLE

doi:10.1038/nature11247

An integrated encyclopedia of DNA elements in the human genome

The ENCODE Project Consortium*

The human genome encodes the blueprint of life, but the function of the vast majority of its nearly three billion bases is unknown. The Encyclopedia of DNA Elements (ENCODE) project has systematically mapped regions of transcription, transcription factor association, chromatin structure and histone modification. These data enabled us to assign biochemical functions for 80% of the genome, in particular outside of the well-studied protein-coding regions. Many discovered candidate regulatory elements are physically associated with one another and with expressed genes, providing new insights into the mechanisms of gene regulation. The newly identified elements also show a statistical correspondence to sequence variants linked to human disease, and can thereby guide interpretation of this variation. Overall, the project provides new insights into the organization and regulation of our genes and genome, and is an expansive resource of functional annotations for biomedical research.

ARTICLE

doi:10.1038/nature11247

An integrated encyclopedia of DNA elements in the human genome

The ENCODE Project Consortium*

ARTICLE

doi:10.1038/nature11232

The accessible chromatin landscape of the human genome

ARTICLE

doi:10.1038/nature11212

An expansive human regulatory lexicon encoded in transcription factor footprints

ARTICLE

doi:10.1038/nature11245

Architecture of the human regulatory network derived from ENCODE data

LETTER

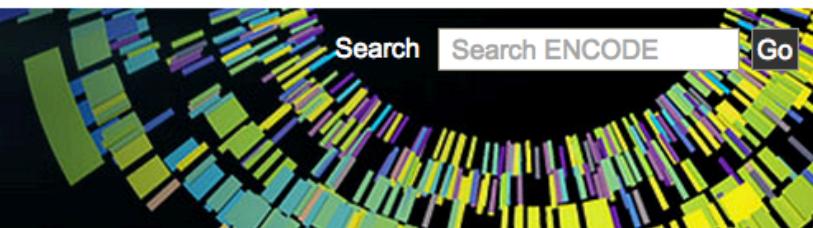
doi:10.1038/nature11279

The long-range interaction landscape of gene promoters

ARTICLE

doi:10.1038/nature11233

Landscape of transcription in human cells

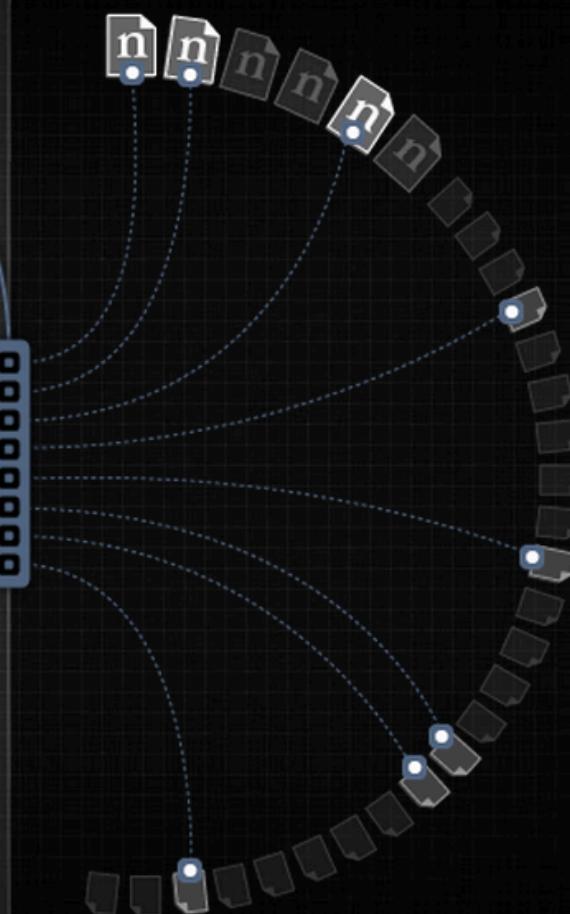


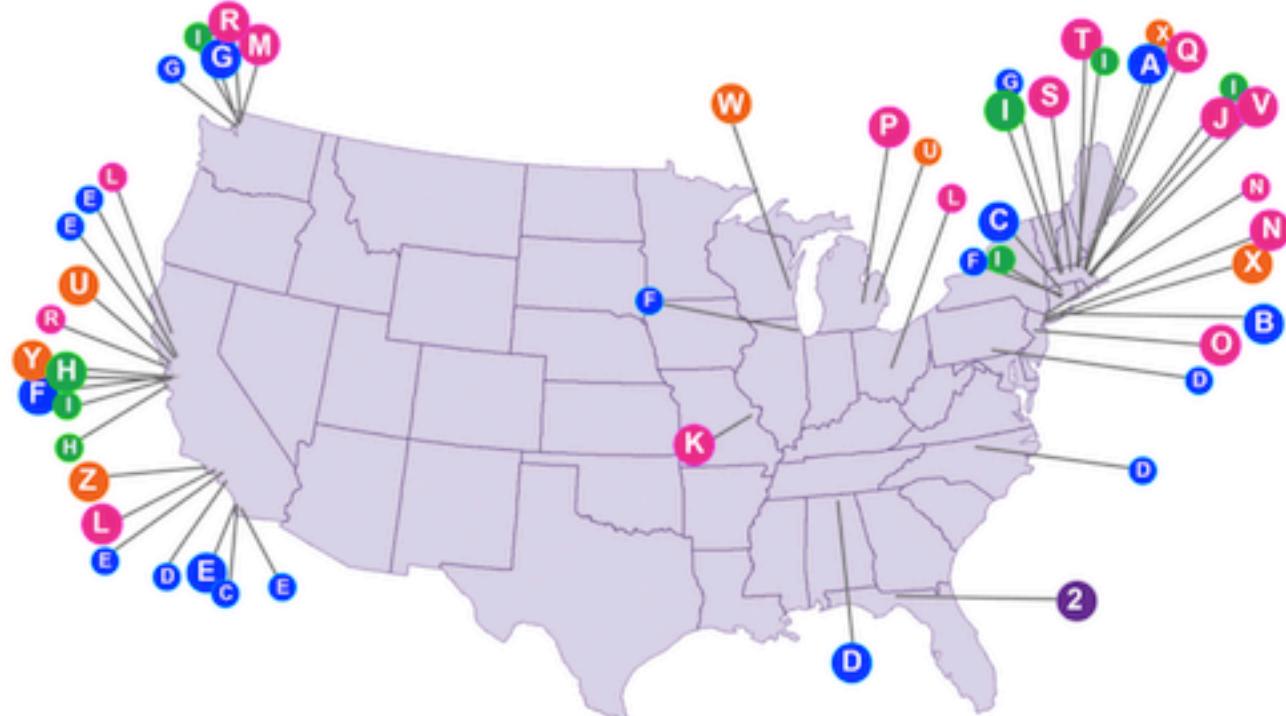
nature ENCODE explorer

THREADS

PAPERS

PRODUCED WITH
SUPPORT FROM
illumina®





Production Groups

- A** Broad Institute
- B** Cold Spring Harbor;
Centre for Genomic Regulation (CRG);
- C** University of Connecticut Health Center;
UCSD
- D** HudsonAlpha; Pennsylvania State;
UC Irvine; Duke; Caltech
- E** UCSD; Salk Institute ; Joint Genome Institute;
Lawrence Berkeley National Laboratory; UCSD
- F** Stanford; University of Chicago; Yale
- G** University of Washington;
Fred Hutchinson Cancer Research Center;
University of Massachusetts Medical School

Data Coordination Center

- H** Stanford; UCSC

Data Analysis Center

- I** University of Massachusetts Medical School;
Yale; MIT; Stanford; Harvard; University of Washington

Technology Development Groups

- J** MIT
- K** Washington University, St. Louis
- L** USC; Ohio State University; UC, Davis
- M** University of Washington
- N** Sloan-Kettering; Weill Cornell Medical College
- O** Princeton; Weizmann
- P** University of Michigan
- Q** Broad Institute
- R** University of Washington; UCSF
- S** Advanced RNA Technologies, LLC
- T** Harvard

Computational Analysis Groups

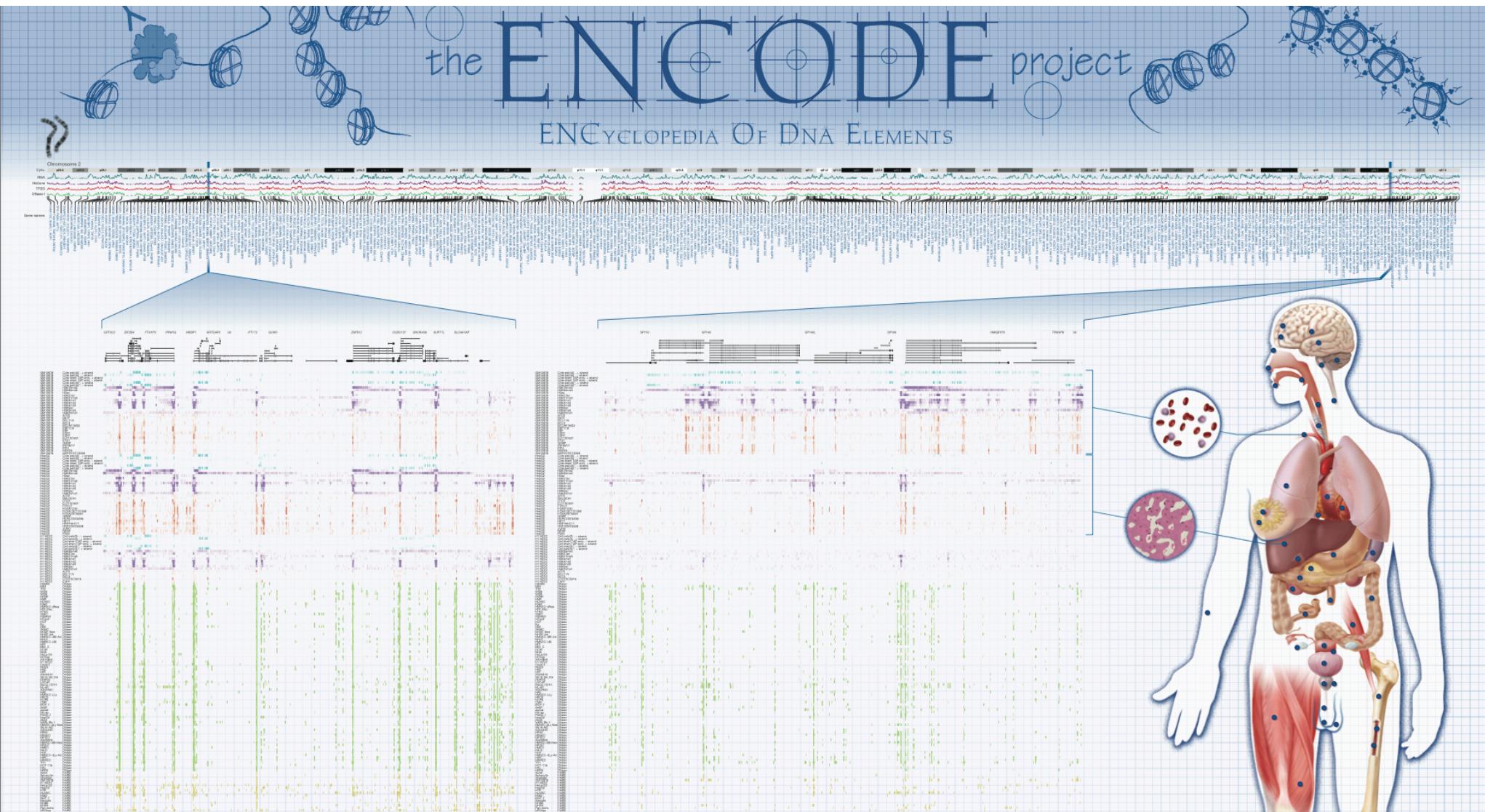
- U** Berkeley; Wayne State University
- V** MIT
- W** University of Wisconsin
- X** Sloan-Kettering; Broad Institute
- Y** Stanford
- Z** UCLA

Affiliated Groups

- 1** Wellcome Trust Sanger Institute
- 2** Florida State University



ENCODE Data Sets



1,640 data sets total over 147 different cell types

ENCODE Data Sets

MAKING A GENOME MANUAL

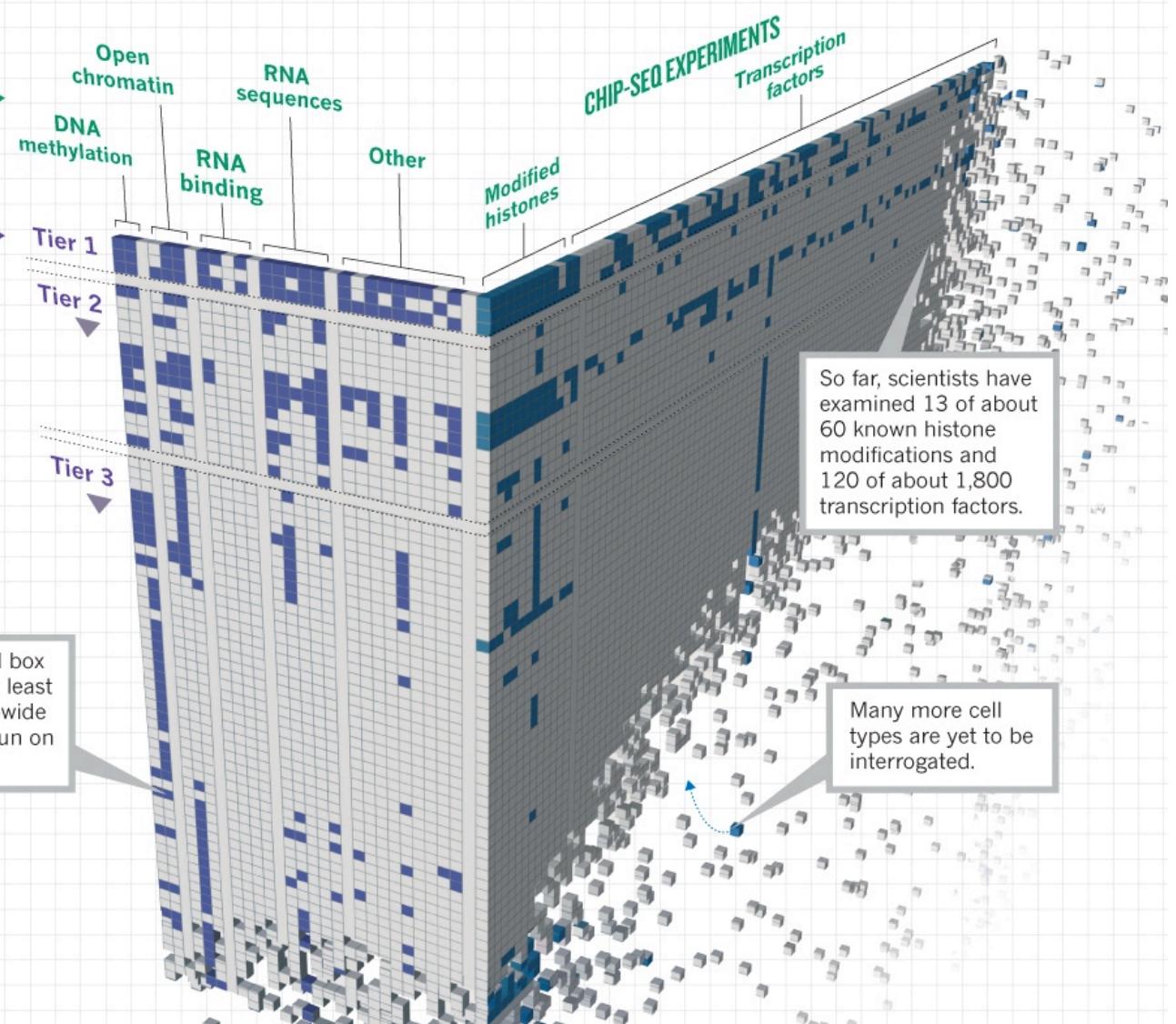
Scientists in the Encyclopedia of DNA Elements Consortium have applied 24 experiment types (across) to more than 150 cell lines (down) to assign functions to as many DNA regions as possible — but the project is still far from complete.

| EXPERIMENTAL TARGETS | |
|-------------------------------|---|
| DNA methylation: | regions layered with chemical methyl groups, which regulate gene expression. |
| Open chromatin: | areas in which the DNA and proteins that make up chromatin are accessible to regulatory proteins. |
| RNA binding: | positions where regulatory proteins attach to RNA. |
| RNA sequences: | regions that are transcribed into RNA. |
| ChIP-seq: | technique that reveals where proteins bind to DNA. |
| Modified histones: | histone proteins, which package DNA into chromosomes, modified by chemical marks. |
| Transcription factors: | proteins that bind to DNA and regulate transcription. |

CELL LINES

- Tiers 1 and 2: widely used cell lines that were given priority.
- Tier 3: all other cell types.

Every shaded box represents at least one genome-wide experiment run on a cell type.



1,640 data sets total over 147 different cell types

Cell Types

Tier 1 (3 samples, most complete analysis)

- **GM12878 (NA12878)**: a lymphoblastoid cell line produced from the blood of a female donor with northern and western European ancestry by EBV transformation. It was one of the original HapMap cell lines and has been deeply sequenced using the Solexa/Illumina platform.
- **K562**: an immortalized cell line produced from a female patient with chronic myelogenous leukemia (CML). It is a widely used model for cell biology, biochemistry, and erythropoiesis. It grows well, is transfectable, and represents the mesoderm lineage.
- **HI-hESC**: HI-human embryonic stem cells

Tier 2 (9 samples, intermediate analysis)

- **HeLa-S3**: cervical carcinoma cells
- **HepG2**: hepatoblastoma cells & model system for metabolism disorders
- **HUVECs**: Primary (non-transformed) human umbilical vein endothelial cells
- Several other major cell lines from cancer and normal tissues

Tier 3 (135 samples, partial analysis)

- Everything else: many major cell lines and body organs

Assays

1. RNA transcribed regions

- RNA-seq: General sequencing of RNA
- CAGE: Identify transcription start sites
- RNA-PET: full length RNA analysis and manual annotation

2. Protein-coding regions

- Mass Spectrometry: Sequencing of proteins

3. Transcription-factor-binding sites

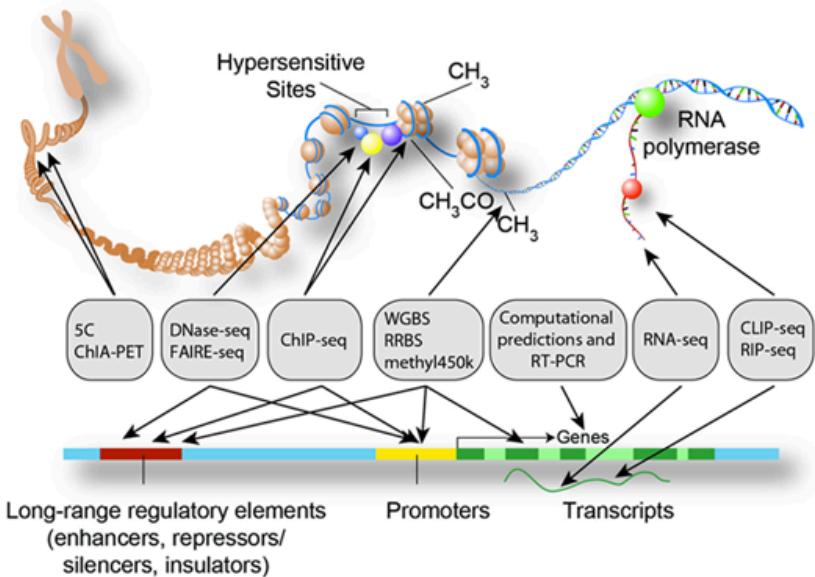
- ChIP-seq: 119 of 1,800 known transcription factors
- DNase-seq: open chromatin accessible to Dnase I cutting, “hallmark of regulatory regions”

4. Chromatin structure

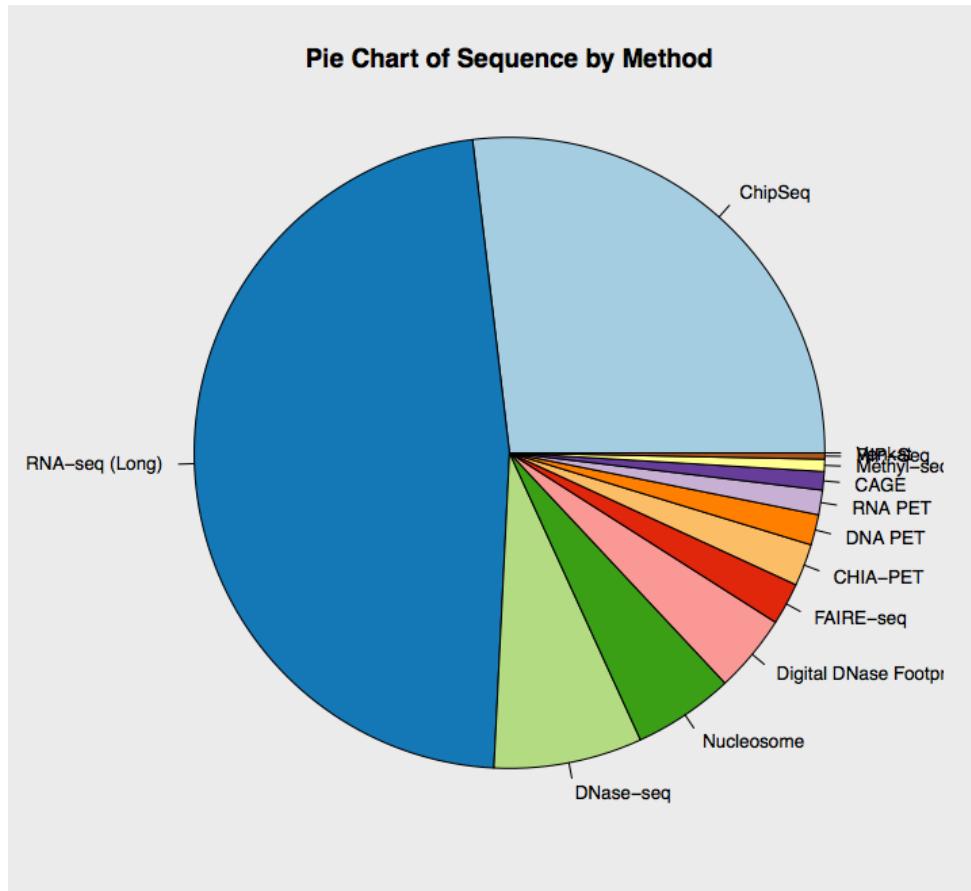
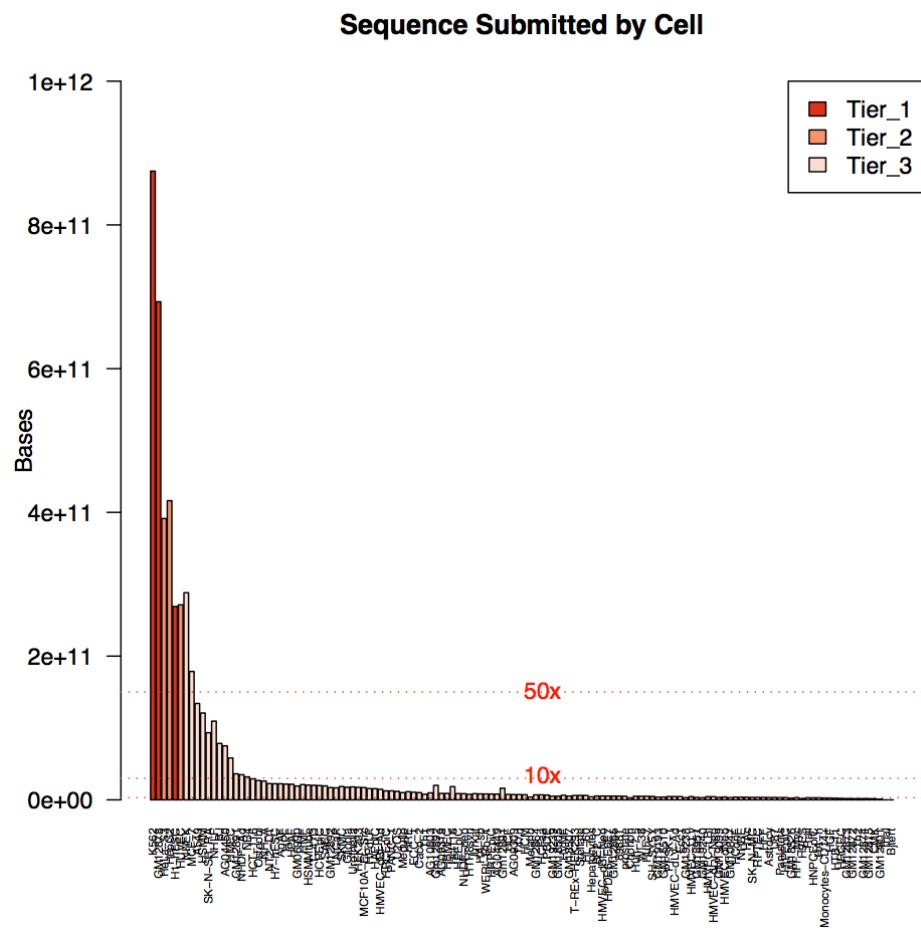
- DNase-seq: 13 of more than 60 currently known histone or DNA modifications
- FAIRE-seq: nucleosome-depleted regions
- Histone ChIP-seq: histone proteins pull down and sequencing
- MNase-seq: nucleosome identification

5. DNA methylation sites

- RRBS assay: Methyl-seq at targeted sites near restriction binding sites



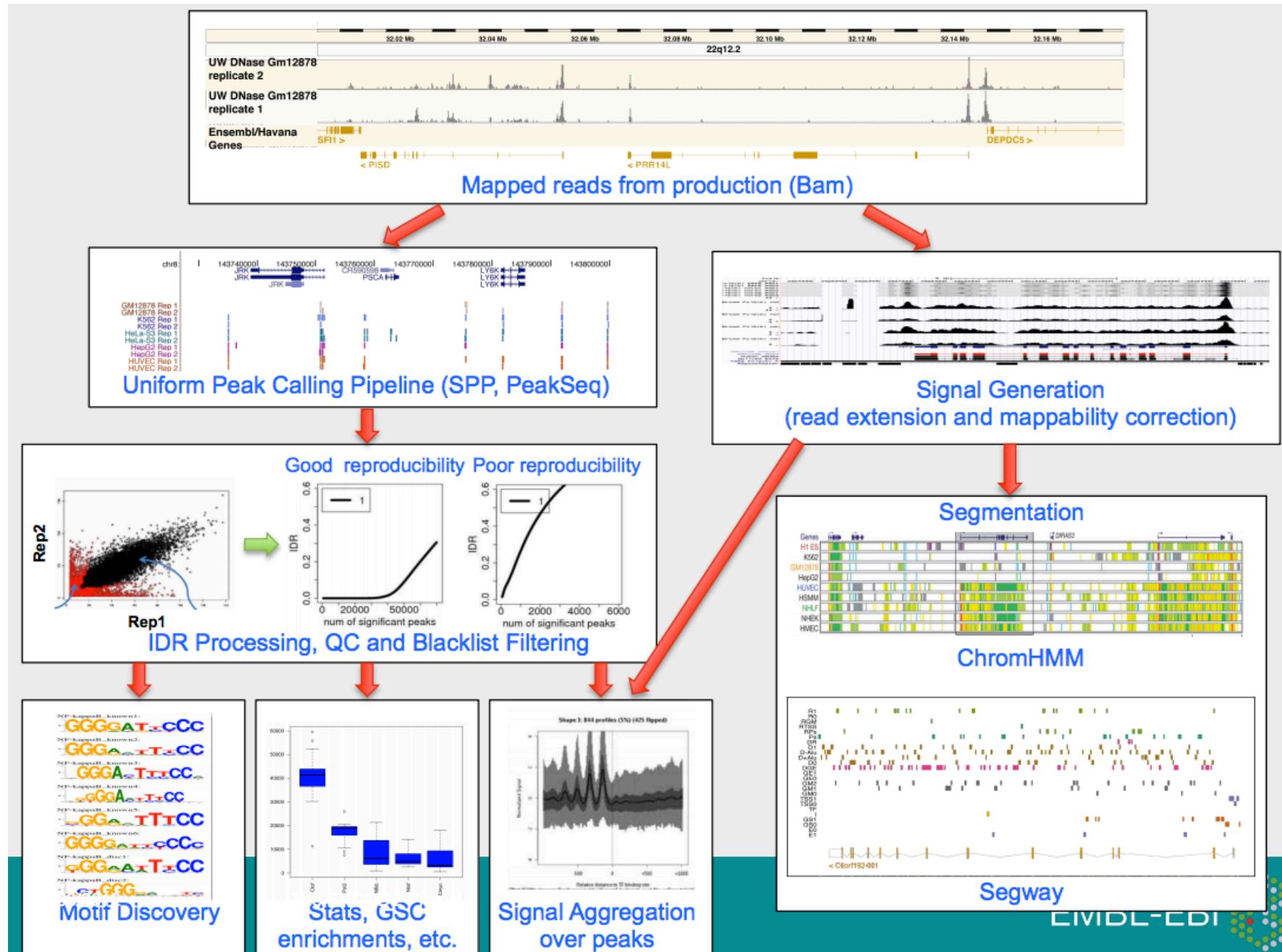
Data Summary



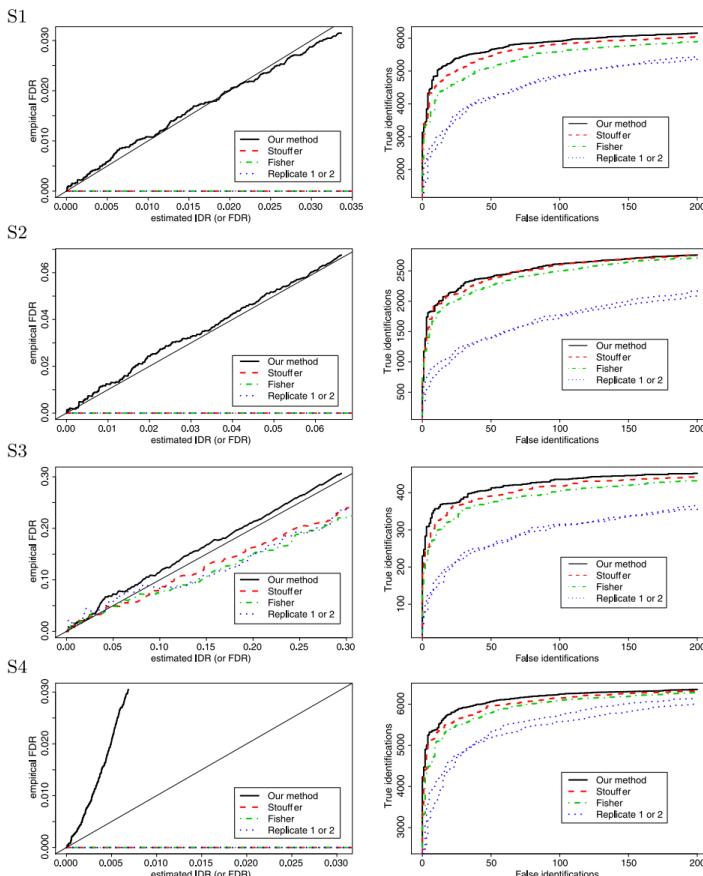
16031 files
1847 Experiments

>5 TeraBases
1716x of the Human Genome

Data Analysis Overview



Irreproducible Discovery Rate (IDR)



Resource

ChIP-seq guidelines and practices of the ENCODE and modENCODE consortia

Stephen G. Landt,^{1,26} Georgi K. Marinov,^{2,26} Anshul Kundaje,^{3,26} Pouya Kheradpour,⁴ Florencia Pauli,⁵ Serafim Batzoglou,³ Bradley E. Bernstein,⁶ Peter Bickel,⁷ James B. Brown,⁷ Philip Cayting,¹ Yiwen Chen,⁸ Gilberto DeSalvo,² Charles Epstein,⁶ Katherine I. Fisher-Aylor,² Ghia Euskirchen,¹ Mark Gerstein,⁹ Jason Gertz,⁵ Alexander J. Hartemink,¹⁰ Michael M. Hoffman,¹¹ Vishwanath R. Iyer,¹² Youngsook L. Jung,^{13,14} Subhradip Karmakar,¹⁵ Manolis Kellis,⁴ Peter V. Kharchenko,¹² Qunhua Li,¹⁶ Tao Liu,⁸ X. Shirley Liu,⁸ Lijia Ma,¹⁵ Aleksandar Milosavljevic,¹⁷ Richard M. Myers,⁵ Peter J. Park,^{13,14} Michael J. Pazin,¹⁸ Marc D. Perry,¹⁹ Debasish Raha,²⁰ Timothy E. Reddy,^{5,27} Joel Rozowsky,⁹ Noam Shores,⁶ Arend Sidow,^{1,21} Matthew Slattery,¹⁵ John A. Stamatoyannopoulos,^{11,22} Michael Y. Tolstorukov,^{13,14} Kevin P. White,¹⁵ Simon Xi,²³ Peggy J. Farnham,^{24,28} Jason D. Lieb,^{25,28} Barbara J. Wold,^{2,28} and Michael Snyder^{1,28}

^{1–25}[Author affiliations appear at the end of the paper.]

Chromatin immunoprecipitation (ChIP) followed by high-throughput DNA sequencing (ChIP-seq) has become a valuable and widely used approach for mapping the genomic location of transcription-factor binding and histone modifications in living cells. Despite its widespread use, there are considerable differences in how these experiments are conducted, how the results are scored and evaluated for quality, and how the data and metadata are archived for public use. These practices affect the quality and utility of any global ChIP experiment. Through our experience in performing ChIP-seq experiments, the ENCODE and modENCODE consortia have developed a set of working standards and guidelines for ChIP experiments that are updated routinely. The current guidelines address antibody validation, experimental replication, sequencing depth, data and metadata reporting, and data quality assessment. We discuss how ChIP quality, assessed in these ways, affects different uses of ChIP-seq data. All data sets used in the analysis have been deposited for public viewing and downloading at the ENCODE (<http://encodeproject.org/ENCODE/>) and modENCODE (<http://www.modencode.org/>) portals.

[Supplemental material is available for this article.]

ARTICLE

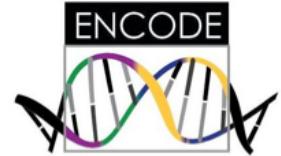
doi:10.1038/nature11247

An integrated encyclopedia of DNA elements in the human genome

The ENCODE Project Consortium*

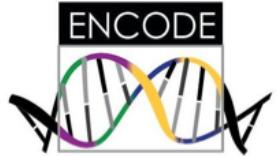
The human genome encodes the blueprint of life, but the function of the vast majority of its nearly three billion bases is unknown. The Encyclopedia of DNA Elements (ENCODE) project has systematically mapped regions of transcription, transcription factor association, chromatin structure and histone modification. These data enabled us to assign biochemical functions for 80% of the genome, in particular outside of the well-studied protein-coding regions. Many discovered candidate regulatory elements are physically associated with one another and with expressed genes, providing new insights into the mechanisms of gene regulation. The newly identified elements also show a statistical correspondence to sequence variants linked to human disease, and can thereby guide interpretation of this variation. Overall, the project provides new insights into the organization and regulation of our genes and genome, and is an expansive resource of functional annotations for biomedical research.

Major Findings



1. *The vast majority (80.4%) of the human genome participates in at least one biochemical RNA- and/or chromatin-associated event in at least one cell type.*
2. *Primate-specific elements as well as elements without detectable mammalian constraint show, in aggregate, evidence of negative selection; thus, some of them are expected to be functional.*
3. *Classifying the genome into seven chromatin states indicates an initial set of 399,124 regions with enhancer-like features and 70,292 regions with promoter-like features, as well as hundreds of thousands of quiescent regions. High-resolution analyses further subdivide the genome into thousands of narrow states with distinct functional properties.*
4. *It is possible to correlate quantitatively RNA sequence production and processing with both chromatin marks and transcription factor binding at promoters, indicating that promoter functionality can explain most of the variation in RNA expression.*
5. *Many non-coding variants in individual genome sequences lie in ENCODE-annotated functional regions; this number is at least as large as those that lie in protein-coding genes.*
6. *Single nucleotide polymorphisms (SNPs) associated with disease by GWAS are enriched within non-coding functional elements, with a majority residing in or near ENCODE-defined regions that are outside of protein-coding genes. In many cases, the disease phenotypes can be associated with a specific cell type or transcription factor.*

Major Findings



1. ***The vast majority (80.4%) of the human genome participates in at least one biochemical RNA- and/or chromatin-associated event in at least one cell type.***
2. *Primate-specific elements as well as elements without detectable mammalian constraint show, in aggregate, evidence of negative selection; thus, some of them are expected to be functional.*
3. *Classifying the genome into seven chromatin states indicates an initial set of 399,124 regions with enhancer-like features and 70,292 regions with promoter-like features, as well as hundreds of thousands of quiescent regions. High-resolution analyses further subdivide the genome into thousands of narrow states with distinct functional properties.*
4. *It is possible to correlate quantitatively RNA sequence production and processing with both chromatin marks and transcription factor binding at promoters, indicating that promoter functionality can explain most of the variation in RNA expression.*
5. *Many non-coding variants in individual genome sequences lie in ENCODE-annotated functional regions; this number is at least as large as those that lie in protein-coding genes.*
6. *Single nucleotide polymorphisms (SNPs) associated with disease by GWAS are enriched within non-coding functional elements, with a majority residing in or near ENCODE-defined regions that are outside of protein-coding genes. In many cases, the disease phenotypes can be associated with a specific cell type or transcription factor.*

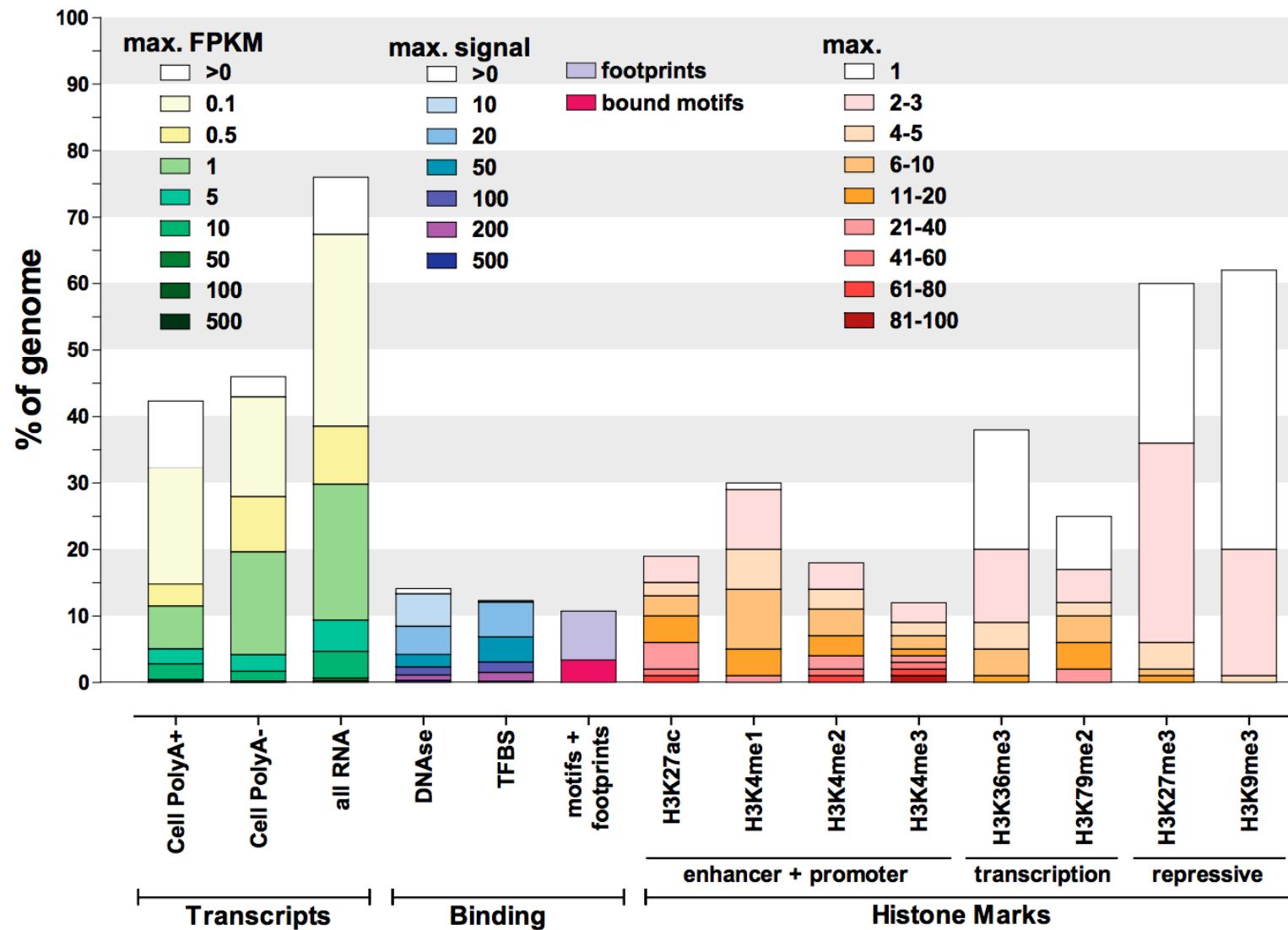
Summary of ENCODE elements

“Accounting for all these elements, a surprisingly large amount of the human genome, 80.4%, is covered by at least one ENCODE-identified element”

- 62% transcribed
- 56% enriched for histone marks
- 15% open chromatin
- 8% TF binding
- 19% At least one DHS or TF Chip-seq peak
- 4% TF binding site motif
- (Note protein coding genes comprise ~2.94% of the genome)

*“Given that the ENCODE project did not assay all cell types, or all transcription factors, and in particular has sampled few specialized or developmentally restricted cell lineages, **these proportions must be underestimates of the total amount of functional bases.**”*

Pervasive Transcription and Regulation



Defining functional DNA elements in the human genome

Kellis et al (2014). PNAS 6131–6138, doi: 10.1073/pnas.1318948111

Redefining the concept of a gene

As a consequence of both the expansion of genic regions by the discovery of new isoforms and the identification of novel intergenic transcripts, there has been a marked increase in the number of intergenic regions (from 32,481 to 60,250) due to their fragmentation and a decrease in their lengths (from 14,170 bp to 3,949 bp median length; Fig. 6). Concordantly, we observed an increased overlap of genic regions. As the determination of genic regions is currently defined by the cumulative lengths of the isoforms and their genetic association to phenotypic characteristics, the likely continued reduction in the lengths of intergenic regions will steadily lead to the overlap of most genes previously assumed to be distinct genetic loci. This supports and is consistent with earlier observations of a highly interleaved transcribed genome¹², but more importantly, prompts the reconsideration of the definition of a gene. As this is a consistent characteristic of annotated genomes, we would propose that the transcript be considered as the basic atomic unit of inheritance. Concomitantly, the term gene would then denote a higher-order concept intended to capture all those transcripts (eventually divorced from their genomic locations) that contribute to a given phenotypic trait. Co-published ENCODE-related papers can be explored online via the *Nature* ENCODE explorer (<http://www.nature.com/ENCODE>), a specially designed visualization tool that allows users to access the linked papers and investigate topics that are discussed in multiple papers via thematically organized threads.

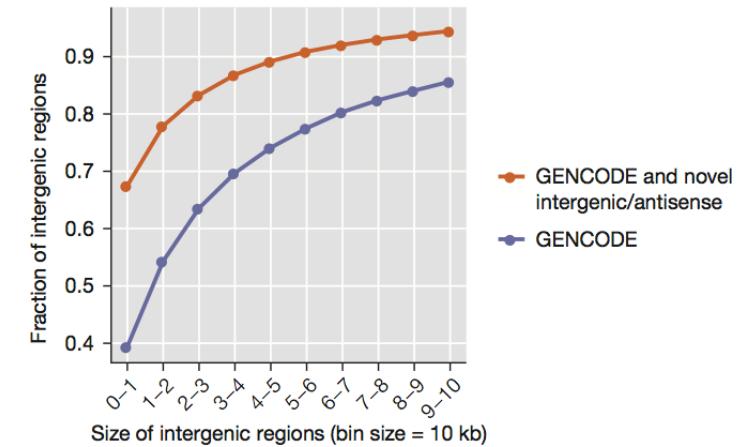


Figure 6 | Size distribution of intergenic regions. Novel genes increase the proportion of small intergenic regions.

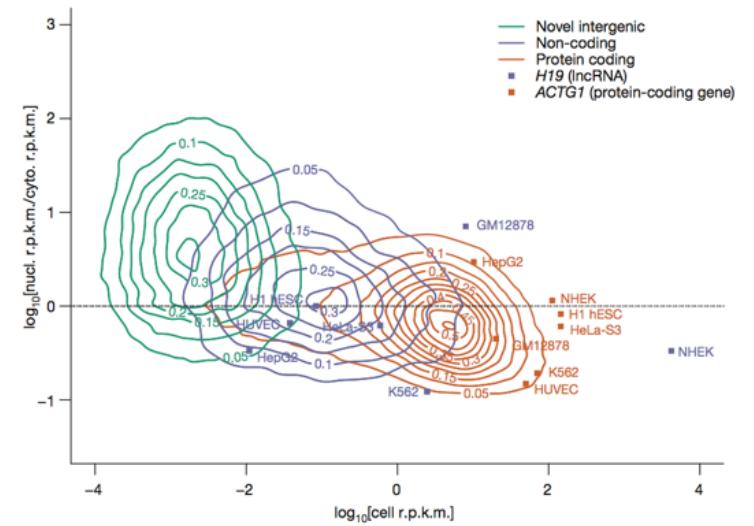
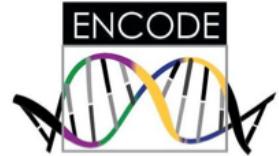


Figure 3 | Abundance of gene types in cellular compartments. Two-dimensional kernel density plots of nuclear over cytosolic enrichment (y axis) versus overall gene expression in the whole cell extract (x axis), for protein coding, long non-coding and novel genes over all cell lines. Only genes present

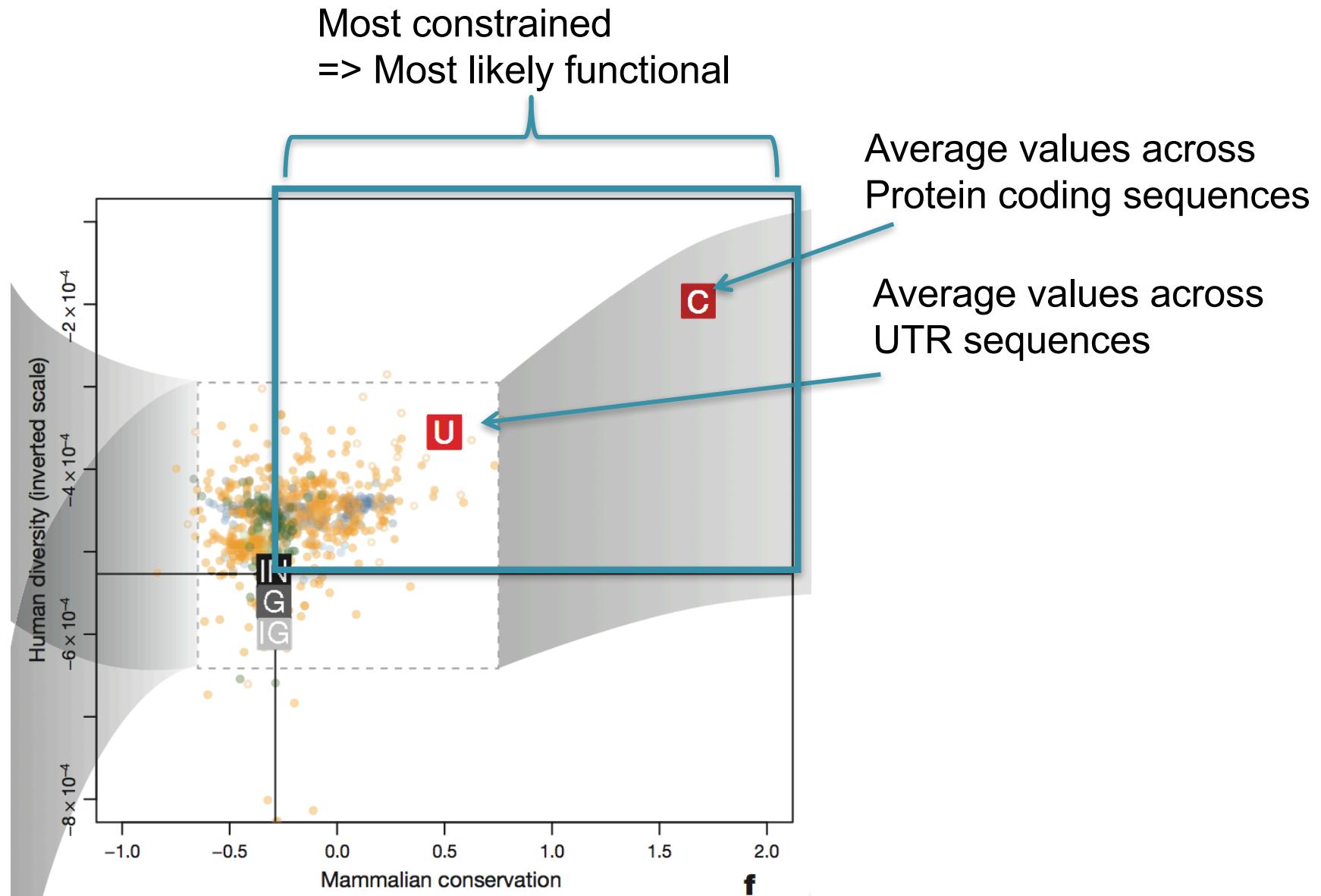
Major Findings



1. *The vast majority (80.4%) of the human genome participates in at least one biochemical RNA- and/or chromatin-associated event in at least one cell type.*
2. ***Primate-specific elements as well as elements without detectable mammalian constraint show, in aggregate, evidence of negative selection; thus, some of them are expected to be functional.***
3. *Classifying the genome into seven chromatin states indicates an initial set of 399,124 regions with enhancer-like features and 70,292 regions with promoter-like features, as well as hundreds of thousands of quiescent regions. High-resolution analyses further subdivide the genome into thousands of narrow states with distinct functional properties.*
4. *It is possible to correlate quantitatively RNA sequence production and processing with both chromatin marks and transcription factor binding at promoters, indicating that promoter functionality can explain most of the variation in RNA expression.*
5. *Many non-coding variants in individual genome sequences lie in ENCODE-annotated functional regions; this number is at least as large as those that lie in protein-coding genes.*
6. *Single nucleotide polymorphisms (SNPs) associated with disease by GWAS are enriched within non-coding functional elements, with a majority residing in or near ENCODE-defined regions that are outside of protein-coding genes. In many cases, the disease phenotypes can be associated with a specific cell type or transcription factor.*

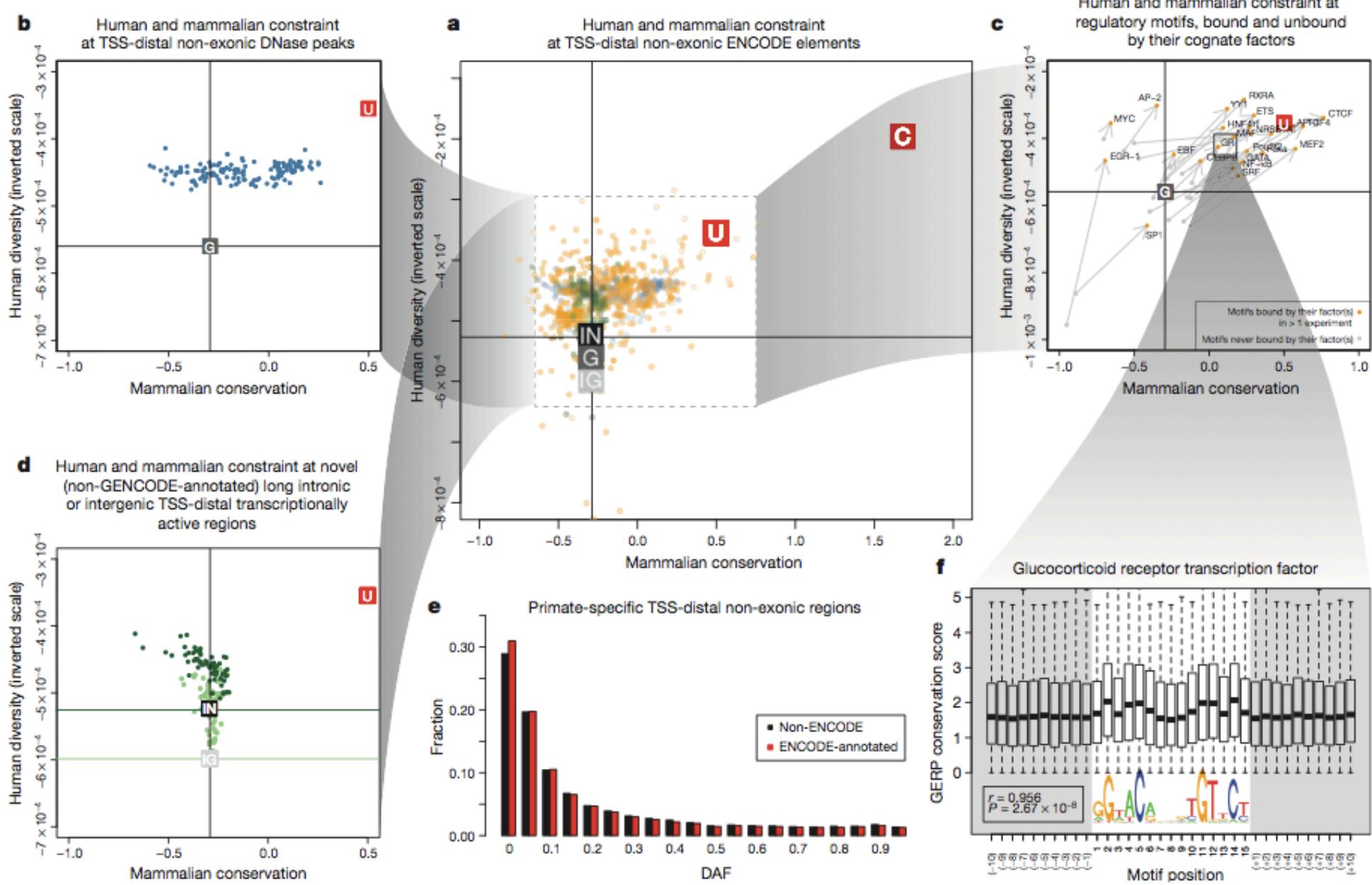
Impact and Evidence of Selection

For a given ENCODE region, how much conservation do we see across modern humans (1000 genomes project)

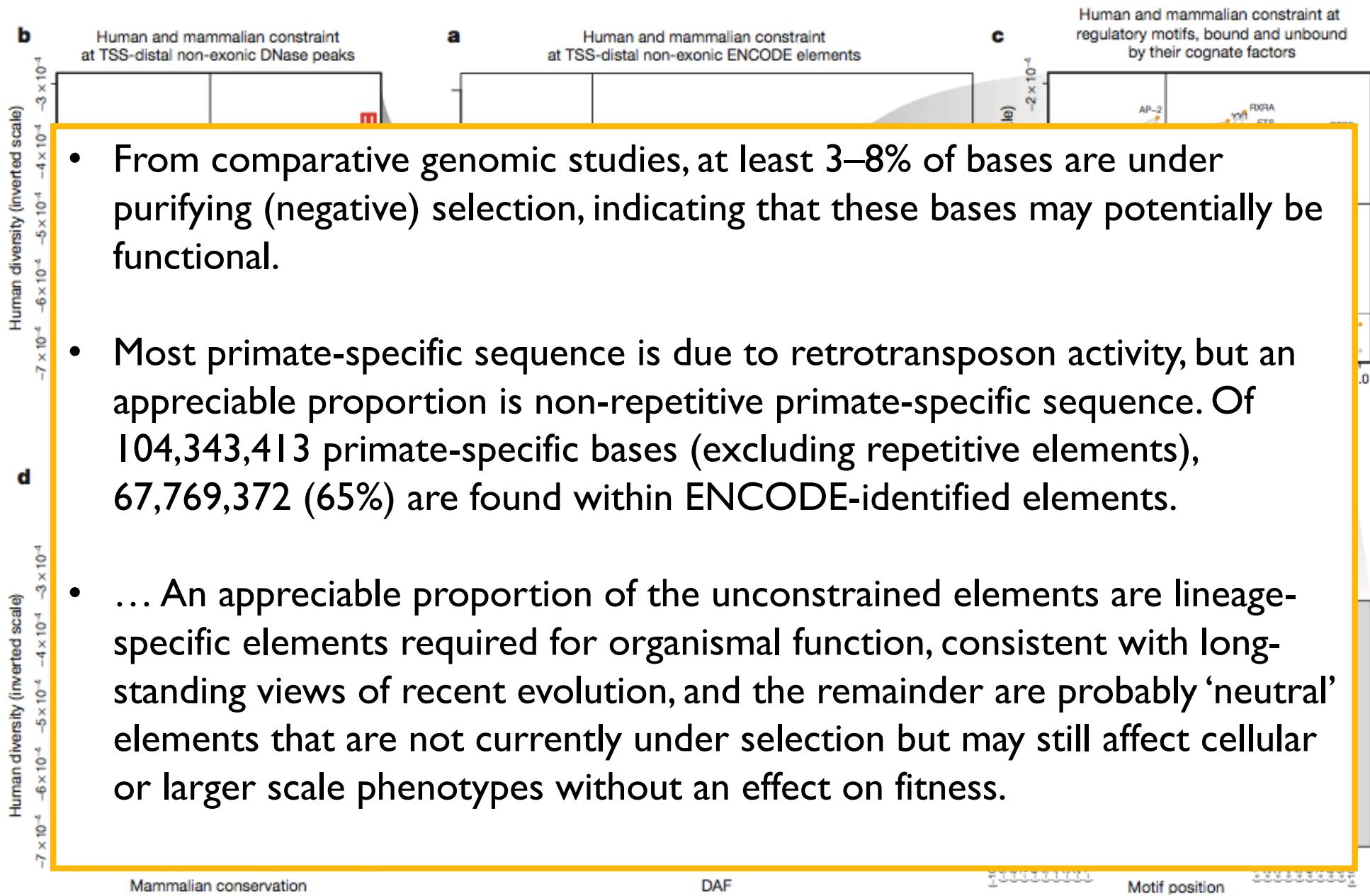


For a given ENCODE region, how much conservation do we see across 24 sequenced mammalian genomes?

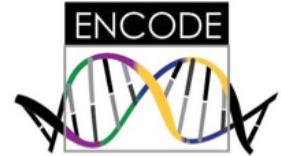
Impact and Evidence of Selection



Impact and Evidence of Selection



Major Findings



1. *The vast majority (80.4%) of the human genome participates in at least one biochemical RNA- and/or chromatin-associated event in at least one cell type.*
2. *Primate-specific elements as well as elements without detectable mammalian constraint show, in aggregate, evidence of negative selection; thus, some of them are expected to be functional.*
3. ***Classifying the genome into seven chromatin states indicates an initial set of 399,124 regions with enhancer-like features and 70,292 regions with promoter-like features, as well as hundreds of thousands of quiescent regions. High-resolution analyses further subdivide the genome into thousands of narrow states with distinct functional properties.***
4. *It is possible to correlate quantitatively RNA sequence production and processing with both chromatin marks and transcription factor binding at promoters, indicating that promoter functionality can explain most of the variation in RNA expression.*
5. *Many non-coding variants in individual genome sequences lie in ENCODE-annotated functional regions; this number is at least as large as those that lie in protein-coding genes.*
6. *Single nucleotide polymorphisms (SNPs) associated with disease by GWAS are enriched within non-coding functional elements, with a majority residing in or near ENCODE-defined regions that are outside of protein-coding genes. In many cases, the disease phenotypes can be associated with a specific cell type or transcription factor.*

Signal Integration

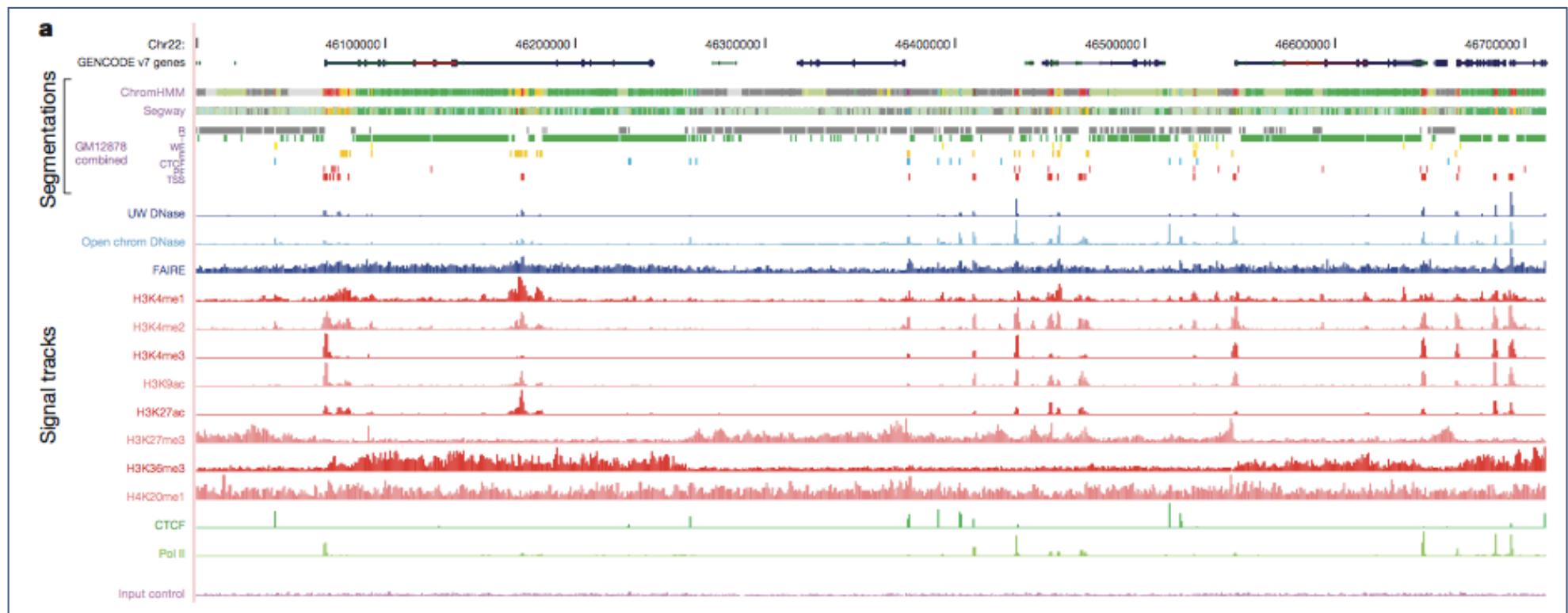
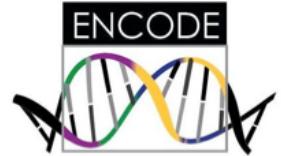


Table 3 | Summary of the combined state types

| Label | Description | Details* | Colour |
|---------|--|---|-------------------|
| CTCF | CTCF-enriched element | Sites of CTCF signal lacking histone modifications, often associated with open chromatin. Many probably have a function in insulator assays, but because of the multifunctional nature of CTCF, we are conservative in our description. Also enriched for the cohesin components RAD21 and SMC3; CTCF is known to recruit the cohesin complex. | Turquoise |
| E | Predicted enhancer | Regions of open chromatin associated with H3K4me1 signal. Enriched for other enhancer-associated marks, including transcription factors known to act at enhancers. In enhancer assays, many of these (>50%) function as enhancers. A more conservative alternative would be <i>cis</i> -regulatory regions. Enriched for sites for the proteins encoded by EP300, FOS, FOSL1, GATA2, HDAC8, JUNB, JUND, NFE2, SMARCA4, SMARCB1, SIRT6 and TAL1 genes in K562 cells. Have nuclear and whole-cell RNA signal, particularly poly(A) ⁺ fraction. | Orange |
| PF R | Predicted promoter flanking region Predicted repressed or low-activity region | Regions that generally surround TSS segments (see below). This is a merged state that includes H3K27me3 polycomb-enriched regions, along with regions that are silent in terms of observed signal for the input assays to the segmentations (low or no signal). They may have other signals (for example, RNA, not in the segmentation input data). Enriched for sites for the proteins encoded by REST and some other factors (for example, proteins encoded by BRF2, CEBPB, MAFK, TRIM28, ZNF274 and SETDB1 genes in K562 cells). | Light red Grey |
| TSS | Predicted promoter region including TSS | Found close to or overlapping GENCODE TSS sites. High precision/recall for TSSs. Enriched for H3K4me3. Sites of open chromatin. Enriched for transcription factors known to act close to promoters and polymerases Pol II and Pol III. Short RNAs are most enriched in these segments. | Bright red |
| T | Predicted transcribed region | Overlap gene bodies with H3K36me3 transcriptional elongation signal. Enriched for phosphorylated form of Pol II signal (elongating polymerase) and poly(A) ⁺ RNA, especially cytoplasmic. | Dark green |
| WE | Predicted weak enhancer or open chromatin <i>cis</i> -regulatory element | Similar to the E state, but weaker signals and weaker enrichments. | Yellow |

- Use ChromHMM and Segway to Summarize the individual assays into 7 functional/regulatory states

Major Findings



1. *The vast majority (80.4%) of the human genome participates in at least one biochemical RNA- and/or chromatin-associated event in at least one cell type.*
2. *Primate-specific elements as well as elements without detectable mammalian constraint show, in aggregate, evidence of negative selection; thus, some of them are expected to be functional.*
3. *Classifying the genome into seven chromatin states indicates an initial set of 399,124 regions with enhancer-like features and 70,292 regions with promoter-like features, as well as hundreds of thousands of quiescent regions. High-resolution analyses further subdivide the genome into thousands of narrow states with distinct functional properties.*
4. ***It is possible to correlate quantitatively RNA sequence production and processing with both chromatin marks and transcription factor binding at promoters, indicating that promoter functionality can explain most of the variation in RNA expression.***
5. *Many non-coding variants in individual genome sequences lie in ENCODE-annotated functional regions; this number is at least as large as those that lie in protein-coding genes.*
6. *Single nucleotide polymorphisms (SNPs) associated with disease by GWAS are enriched within non-coding functional elements, with a majority residing in or near ENCODE-defined regions that are outside of protein-coding genes. In many cases, the disease phenotypes can be associated with a specific cell type or transcription factor.*

Histone Modifications

- Histones are the proteins around which DNA is wound into nucleosomes and at a higher level chromatin
- Histone modifications have been previously reported to indicate repressive/activating functional state
- Use ChIP-seq techniques to locate where they are in the genome
 - Cannot be predicted from sequence composition alone, highly dependent on cell type and cell state)

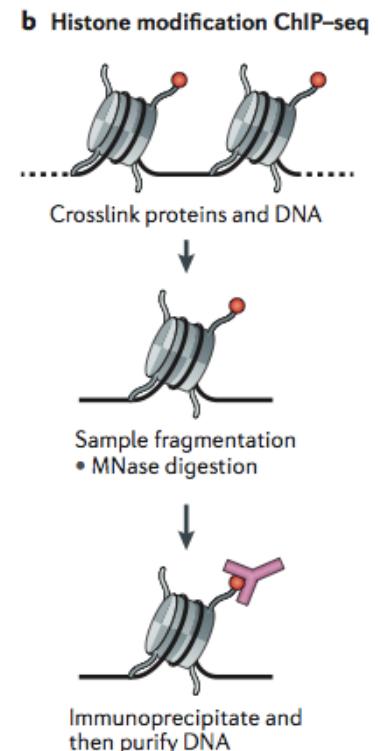
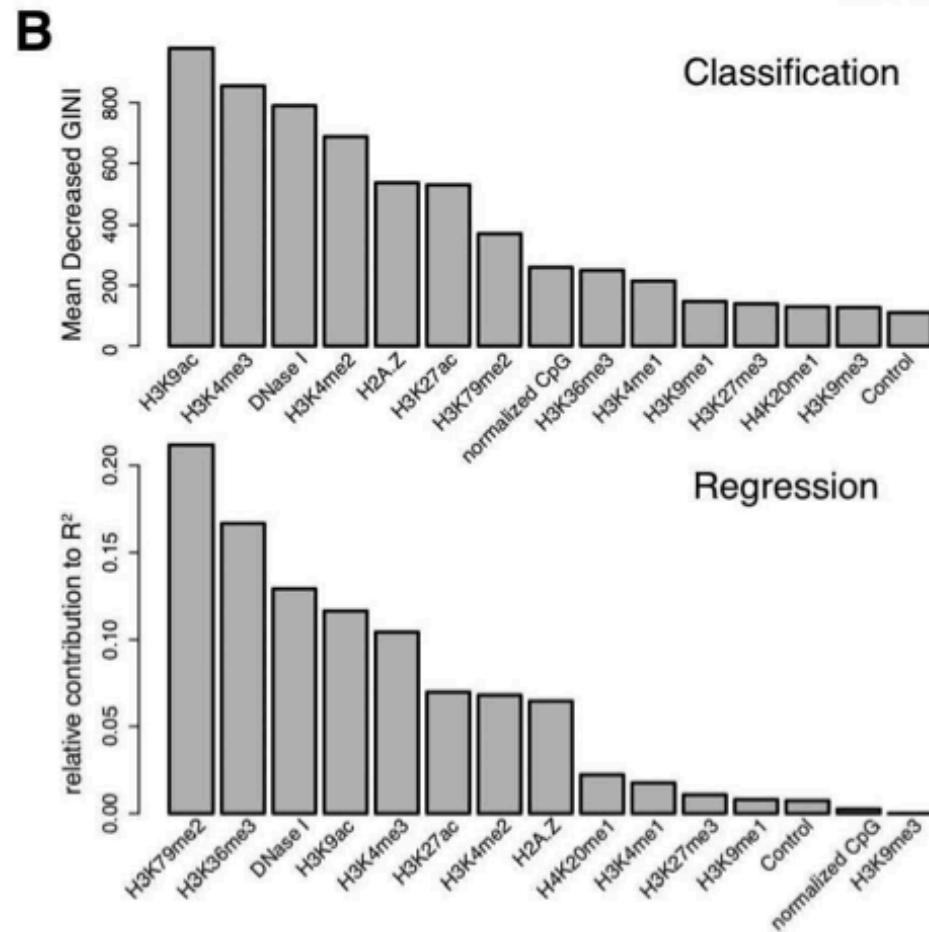
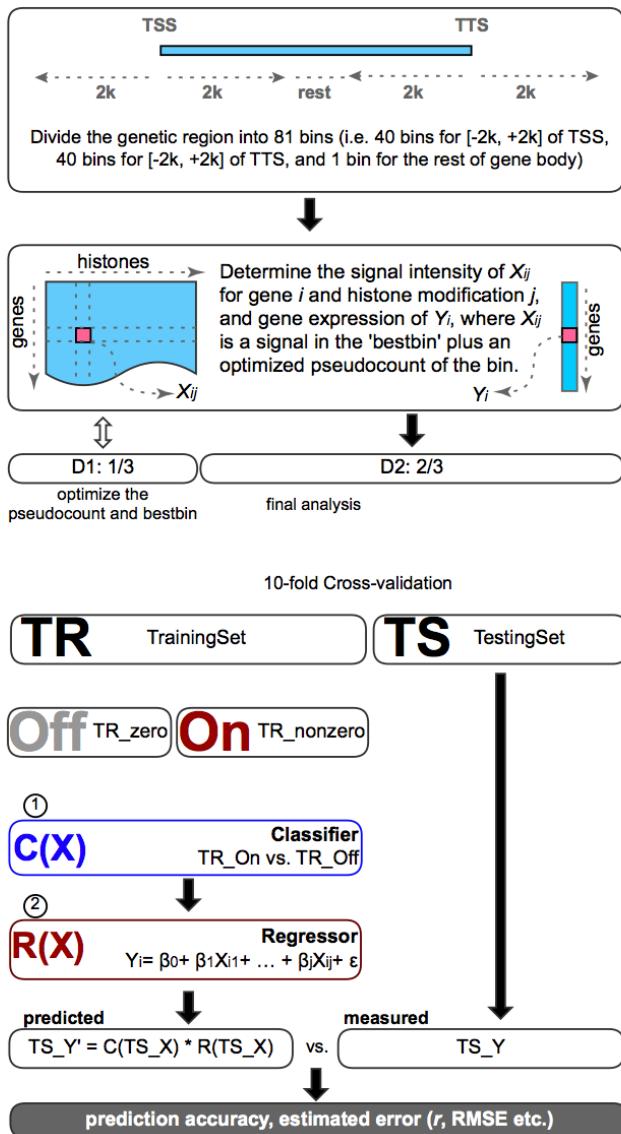


Table 2 | Summary of ENCODE histone modifications and variants

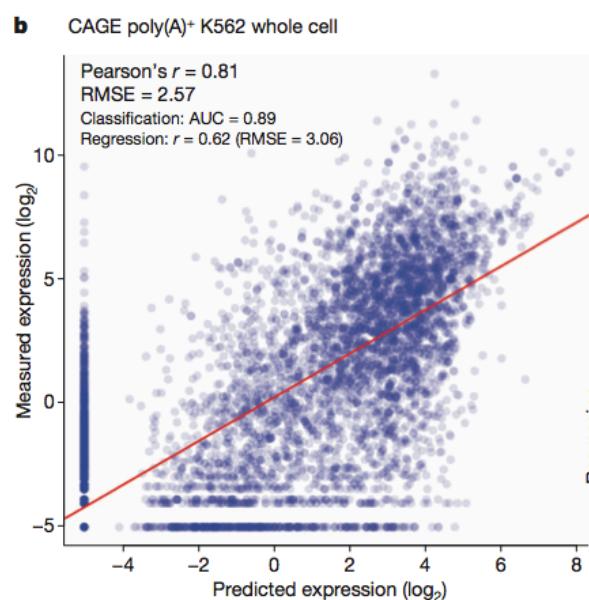
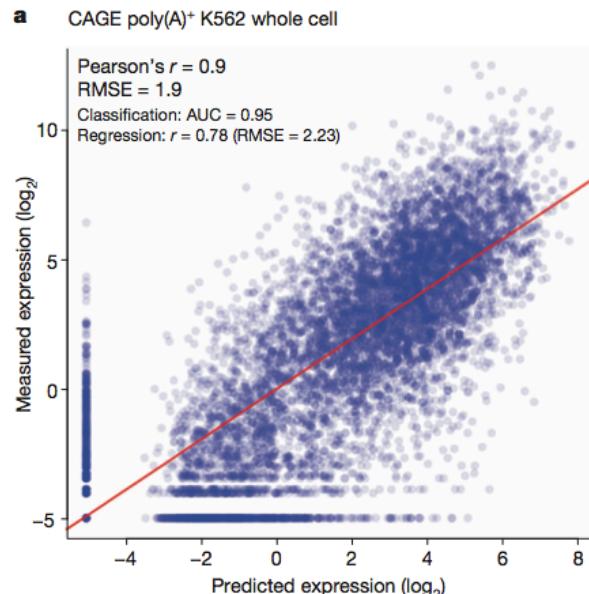
| Histone modification or variant | Signal characteristics | Putative functions |
|---------------------------------|------------------------|---|
| H2A.Z | Peak | Histone protein variant (H2A.Z) associated with regulatory elements with dynamic chromatin |
| H3K4me1 | Peak/region | Mark of regulatory elements associated with enhancers and other distal elements, but also enriched downstream of transcription starts |
| H3K4me2 | Peak | Mark of regulatory elements associated with promoters and enhancers |
| H3K4me3 | Peak | Mark of regulatory elements primarily associated with promoters/transcription starts |
| H3K9ac | Peak | Mark of active regulatory elements with preference for promoters |
| H3K9me1 | Region | Preference for the 5' end of genes |
| H3K9me3 | Peak/region | Repressive mark associated with constitutive heterochromatin and repetitive elements |
| H3K27ac | Peak | Mark of active regulatory elements; may distinguish active enhancers and promoters from their inactive counterparts |
| H3K27me3 | Region | Repressive mark established by polycomb complex activity associated with repressive domains and silent developmental genes |
| H3K36me3 | Region | Elongation mark associated with transcribed portions of genes, with preference for 3' regions after intron 1 |
| H3K79me2 | Region | Transcription-associated mark, with preference for 5' end of genes |
| H4K20me1 | Region | Preference for 5' end of genes |

Expression Modeling



Modeling gene expression using chromatin features in various cellular context
 Dong et al. (2012) Genome Biology. 12:R53

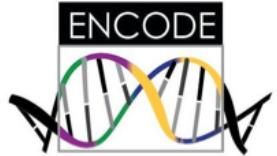
Expression Modeling



- Developed predictive models to explore the interaction between histone modifications and transcription factor binding towards level of transcription
- The best models had two components: an initial classification component (on/off) and a second quantitative model component
- Together, these correlation models indicate both that a limited set of chromatin marks are sufficient to ‘explain’ transcription and that a variety of transcription factors might have broad roles in general transcription levels across many genes

Figure 2 | Modelling transcription levels from histone modification and transcription-factor-binding patterns. a, b, Correlative models between either histone modifications or transcription factors, respectively, and RNA production as measured by CAGE tag density at TSSs in K562 cells. In each case the scatter plot shows the output of the correlation models (x axis) compared to observed values (y axis). The bar graphs show the most important histone modifications (a) or transcription factors (b) in both the initial classification phase (top bar graph) or the quantitative regression phase (bottom bar graph), with larger values indicating increasing importance of the variable in the model. Further analysis of other cell lines and RNA measurement types is reported elsewhere^{59,79}. AUC, area under curve; Gini, Gini coefficient; RMSE, root mean square error.

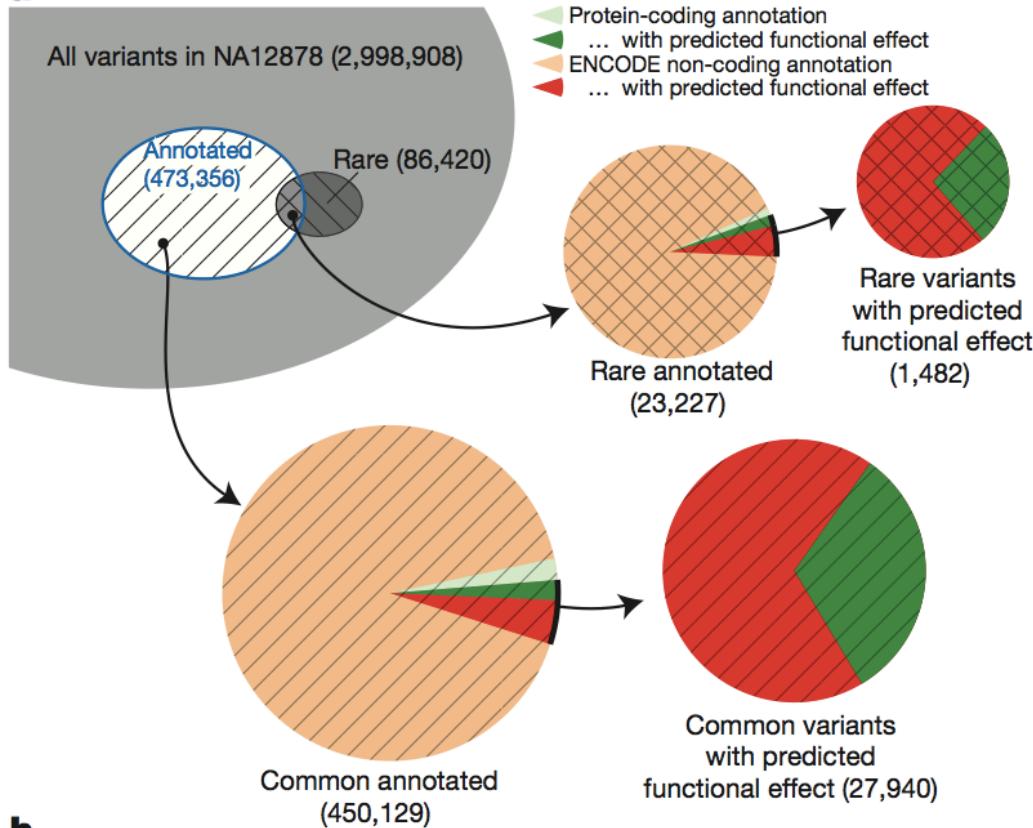
Major Findings



1. *The vast majority (80.4%) of the human genome participates in at least one biochemical RNA- and/or chromatin-associated event in at least one cell type.*
2. *Primate-specific elements as well as elements without detectable mammalian constraint show, in aggregate, evidence of negative selection; thus, some of them are expected to be functional.*
3. *Classifying the genome into seven chromatin states indicates an initial set of 399,124 regions with enhancer-like features and 70,292 regions with promoter-like features, as well as hundreds of thousands of quiescent regions. High-resolution analyses further subdivide the genome into thousands of narrow states with distinct functional properties.*
4. *It is possible to correlate quantitatively RNA sequence production and processing with both chromatin marks and transcription factor binding at promoters, indicating that promoter functionality can explain most of the variation in RNA expression.*
5. ***Many non-coding variants in individual genome sequences lie in ENCODE-annotated functional regions; this number is at least as large as those that lie in protein-coding genes.***
6. *Single nucleotide polymorphisms (SNPs) associated with disease by GWAS are enriched within non-coding functional elements, with a majority residing in or near ENCODE-defined regions that are outside of protein-coding genes. In many cases, the disease phenotypes can be associated with a specific cell type or transcription factor.*

Many variants in ENCODE-regions

a



b

Figure 9 | Examining ENCODE elements on a per individual basis in the normal and cancer genome. a, Breakdown of variants in a single genome (NA12878) by both frequency (common or rare (that is, variants not present in the low-coverage sequencing of 179 individuals in the pilot 1 European panel of the 1000 Genomes project⁵⁵) and by ENCODE annotation, including protein-coding gene and non-coding elements (GENCODE annotations for protein-coding genes, pseudogenes and other ncRNAs, as well as transcription-factor-binding sites from ChIP-seq data sets, excluding broad annotations such as histone modifications, segmentations and RNA-seq). Annotation status is further subdivided by predicted functional effect, being non-synonymous and missense mutations for protein-coding regions and variants overlapping bound transcription factor motifs for non-coding element annotations. A substantial proportion of variants are annotated as having predicted functional effects in the non-coding category. b, One of several relatively rare occurrences, where

Breakdown of variants by frequency

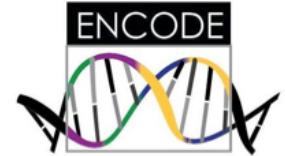
- Common or Rare (that is, variants not present in the low-coverage sequencing of 179 individuals in the pilot 1 European panel of the 1000 Genomes project)
- ENCODE annotation, including protein-coding gene and non-coding elements

Annotation status is further subdivided by predicted functional effect

- non-synonymous and missense mutations for protein-coding regions and variants overlapping bound transcription factor motifs for non-coding element annotations.

A substantial proportion of variants are annotated as having predicted functional effects in the non-coding category.

Major Findings



1. *The vast majority (80.4%) of the human genome participates in at least one biochemical RNA- and/or chromatin-associated event in at least one cell type.*
2. *Primate-specific elements as well as elements without detectable mammalian constraint show, in aggregate, evidence of negative selection; thus, some of them are expected to be functional.*
3. *Classifying the genome into seven chromatin states indicates an initial set of 399,124 regions with enhancer-like features and 70,292 regions with promoter-like features, as well as hundreds of thousands of quiescent regions. High-resolution analyses further subdivide the genome into thousands of narrow states with distinct functional properties.*
4. *It is possible to correlate quantitatively RNA sequence production and processing with both chromatin marks and transcription factor binding at promoters, indicating that promoter functionality can explain most of the variation in RNA expression.*
5. *Many non-coding variants in individual genome sequences lie in ENCODE-annotated functional regions; this number is at least as large as those that lie in protein-coding genes.*
6. ***Single nucleotide polymorphisms (SNPs) associated with disease by GWAS are enriched within non-coding functional elements, with a majority residing in or near ENCODE-defined regions that are outside of protein-coding genes. In many cases, the disease phenotypes can be associated with a specific cell type or transcription factor.***

ENCODE Studies

The screenshot shows the ENCODE website interface. At the top, the title "ENCODE: Encyclopedia of DNA Elements" is displayed. Below the title is a detailed diagram illustrating the ENCODE project's approach to mapping functional elements in the genome. The diagram shows a chromosome with various regulatory elements like Hypersensitive Sites, Long-range regulatory elements (enhancers, repressors/silencers, insulators), Promoters, Genes, and Transcripts. Arrows point from seven different assay categories to these genomic features: 5C ChIA-PET Hi-C, DNase-seq FAIRE-seq ATAC-seq, ChIP-seq, WGBS RRBS methyl array, Computational predictions, RNA-seq, and CLIP-seq RIP-seq. A legend at the bottom indicates the four organisms supported: HUMAN, MOUSE, WORM, and FLY. Below the diagram, there are three circular charts: "Project" (13223), "Biosample Type" (13223), and "Assay Categories" (yellow bar chart showing values from 4000 to 8000). A "View Assay Matrix" button is located between the Project and Biosample Type charts. A "Get Started" button is also visible.

ENCODE: Encyclopedia of DNA Elements

The ENCODE (Encyclopedia of DNA Elements) Consortium is an international collaboration of research groups funded by the National Human Genome Research Institute (NHGRI). The goal of ENCODE is to build a comprehensive parts list of functional elements in the human genome, including elements that act at the protein and RNA levels, and regulatory elements that control cells and circumstances in which a gene is active.

Get Started

Based on an image by Darryl Leja (NHGRI), Ian Dunham (EBI), Michael Pazin (NHGRI)

HUMAN MOUSE WORM FLY

View Assay Matrix

Project

Biosample Type

Assay Categories

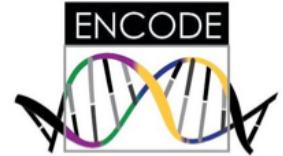
13223

13223

8000
7000
6000
5000
4000

>5000 Citations

Summary & Critique



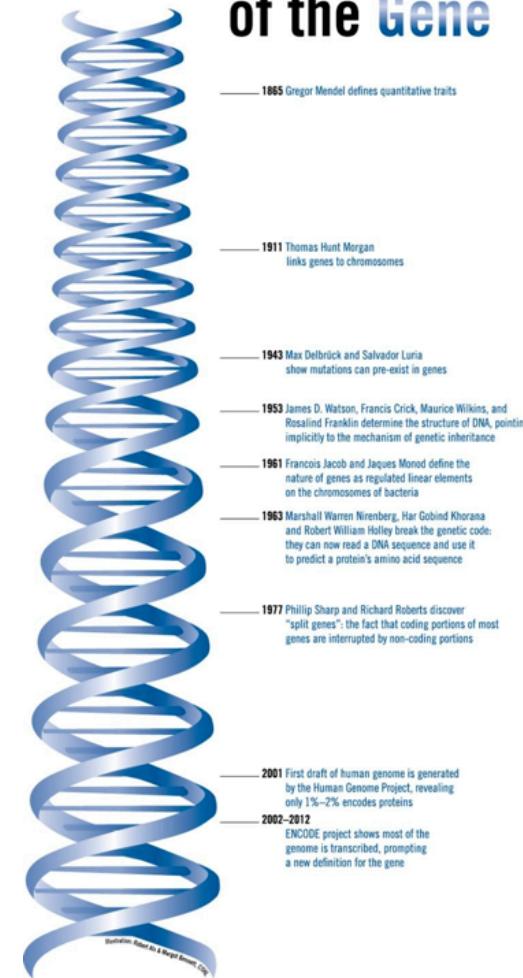
- **Summary**

- *The unprecedented number of functional elements identified in this study provides a valuable resource to the scientific community as well as significantly enhances our understanding of the human genome.*

- **Critique**

- Was it correct?
- What is functional?
- What is conservation?
- What was the control?
- What are the tradeoffs of organizing so much funding (\$288M!) around a single project; will other groups successfully use these data?

Redefining the Nature of the Gene



Comment on “Evidence of Abundant Purifying Selection in Humans for Recently Acquired Regulatory Functions”



Phil Green* and Brent Ewing

Ward and Kellis (Reports, 28 September 2012, p. 1675; published online 5 September 2012) found altered patterns of human polymorphism in biochemically active but non-mammalian-conserved genomic regions relative to control regions and interpreted this as due to lineage-specific purifying selection. We find on closer inspection of their data that the polymorphism trends are primarily attributable to mutational variation and technical artifacts rather than selection.

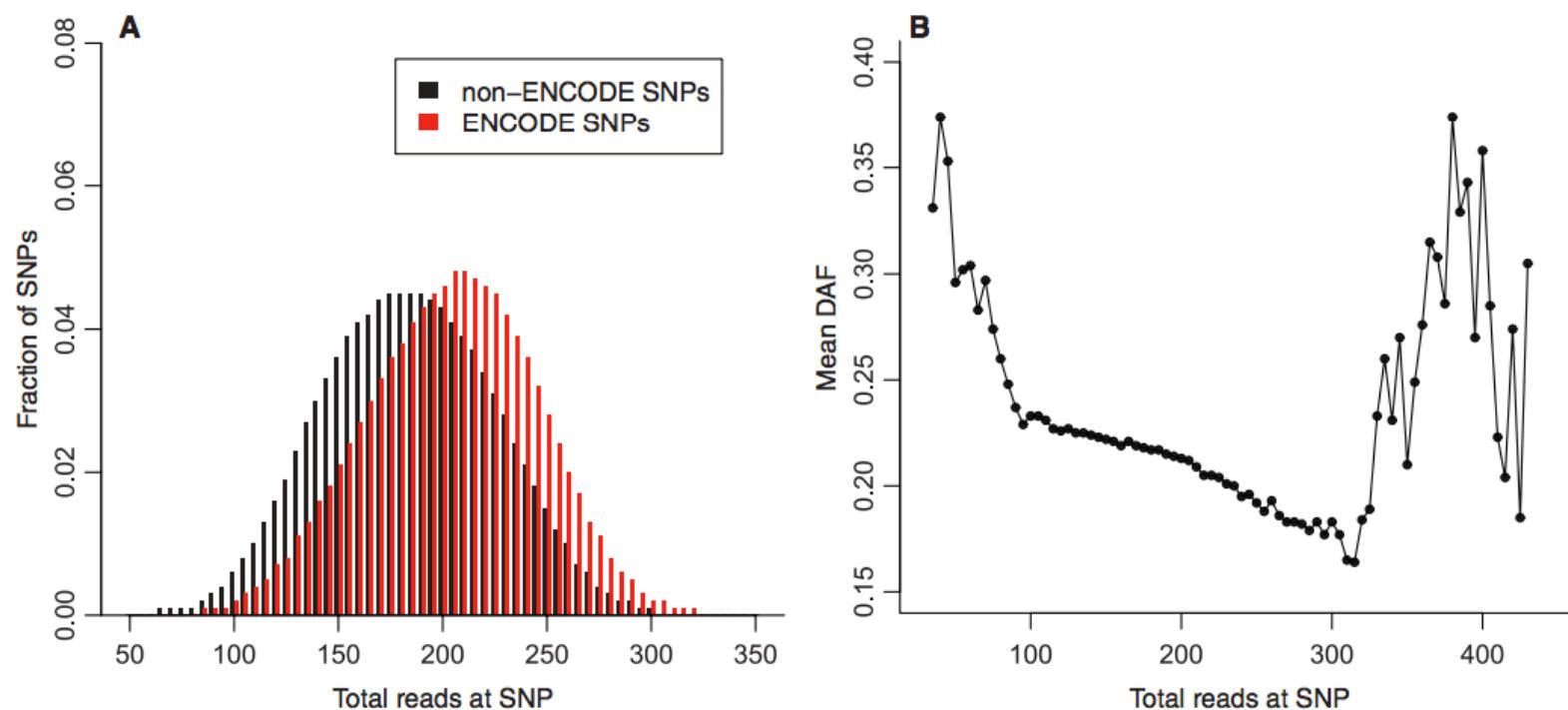


Fig. 1. Variation in 1000 Genomes read depth (totaled over 59 Yoruban individuals) and its impact on DAF. (A) Read-depth distribution for SNPs in neutral control (non-ENCODE) and ENCODE target regions. (B) DAF as a function

of read depth, for non-ENCODE SNPs. DAF decreases with increasing depth, due to increasing sensitivity to detect rare variants; the reverse trend at depths above 300 likely reflects the presence of spurious “paralogue-collapse” SNPs.

On the Immortality of Television Sets: “Function” in the Human Genome According to the Evolution-Free Gospel of ENCODE

Dan Graur^{1,*}, Yichen Zheng¹, Nicholas Price¹, Ricardo B.R. Azevedo¹, Rebecca A. Zufall¹, and Eran Elhaik²

¹Department of Biology and Biochemistry, University of Houston

²Department of Mental Health, Johns Hopkins University Bloomberg School of Public Health

*Corresponding author: E-mail: dgraur@uh.edu.

Accepted: February 16, 2013

Abstract

A recent slew of ENCYclopedia Of DNA Elements (ENCODE) Consortium publications, specifically the article signed by all Consortium members, put forward the idea that more than 80% of the human genome is functional. This claim flies in the face of current estimates according to which the fraction of the genome that is evolutionarily conserved through purifying selection is less than 10%. Thus, according to the ENCODE Consortium, a biological function can be maintained indefinitely without selection, which implies that at least $80 - 10 = 70\%$ of the genome is perfectly invulnerable to deleterious mutations, either because no mutation can ever occur in these “functional” regions or because no mutation in these regions can ever be deleterious. This absurd conclusion was reached through various means, chiefly by employing the seldom used “causal role” definition of biological function and then applying it inconsistently to different biochemical properties, by committing a logical fallacy known as “affirming the consequent,” by failing to appreciate the crucial difference between “junk DNA” and “garbage DNA,” by using analytical methods that yield biased errors and inflate estimates of functionality, by favoring statistical sensitivity over specificity, and by emphasizing statistical significance rather than the magnitude of the effect. Here, we detail the many logical and methodological transgressions involved in assigning functionality to almost every nucleotide in the human genome. The ENCODE results were predicted by one of its authors to necessitate the rewriting of textbooks. We agree, many textbooks dealing with marketing, mass-media hype, and public relations may well have to be rewritten.

Key words: junk DNA, genome functionality, selection, ENCODE project.

Deletion of Ultraconserved Elements Yields Viable Mice

Nadav Ahituv^{1,2✉}, Yiwen Zhu¹, Axel Visel¹, Amy Holt¹, Veena Afzal¹, Len A. Pennacchio^{1,2}, Edward M. Rubin^{1,2*}

1 Genomics Division, Lawrence Berkeley National Laboratory, Berkeley, California, United States of America, **2** United States Department of Energy Joint Genome Institute, Walnut Creek, California, United States of America

Ultraconserved elements have been suggested to retain extended perfect sequence identity between the human, mouse, and rat genomes due to essential functional properties. To investigate the necessities of these elements in vivo, we removed four noncoding ultraconserved elements (ranging in length from 222 to 731 base pairs) from the mouse genome. To maximize the likelihood of identifying functional elements, we chose to remove elements that were active in enhancers in a mouse transgenic assay. We found that these elements were inactivated in the mouse and when they were removed, the resulting lines of mice lacking these elements showed no apparent abnormalities when assayed for a variety of phenotypic traits. In addition, more targeted screens, informed by the results of the initial study, indicated that the remaining investigated elements had been altered without significantly impacting their function. These findings indicate that extreme levels of DNA sequence conservation are not necessarily indicative of an indispensable functional nature.

Citation: Ahituv N, Zhu Y, Visel A, Holt A, Afzal V, et al.

Author Summary

It is widely believed that the most evolutionarily conserved DNA sequences in the human genome have been preserved because of their functional importance and that their removal would thus have a devastating effect on the organism. To ascertain this we removed from the mouse genome four ultraconserved elements—sequences of 200 base pairs or longer that are 100% identical among human, mouse, and rat. To our surprise, we found that the mice lacking these elements are viable, fertile, and show no apparent abnormalities. This completely unexpected finding indicates that extreme levels of DNA sequence conservation are not necessarily indicative of an indispensable functional nature.

The ENCODE project: Missteps overshadowing a success

Two clichés of science journalism have now played out around the ENCODE project. ENCODE's publicity first presented a misleading "all the textbooks"

*"To clarify what noise means, I propose the **Random Genome Project**. Suppose we put a few million bases of entirely random synthetic DNA into a human cell, and do an ENCODE project on it. Will it be reproducibly transcribed into mRNA-like transcripts, reproducibly bound by DNA-binding proteins, and reproducibly wrapped around histones marked by specific chromatin modifications? I think yes.*

A striking feature of genetic regulation is that regulatory factors (proteins or RNAs) generally recognize and bind to small sites, small enough that any given factor will find specific binding sites even in random DNA. Promoters, enhancers, splice sites, poly-A addition sites, and other functional features in the genome all have substantial random occurrence frequencies. These sites are not nonspecific in a random genome. They are specific sequences, albeit randomly occurring and not under selection for any function.

Would biochemical activities in the random genome be regulated under different conditions? For example, would they be cell type-specific? Surely yes, because the regulatory factors themselves (such as transcription factors) are regulated and expressed in specific cell types and conditions."

The ENCODE project: Missteps overshadowing a success

"There are three categories of big science: the big experiment, the map, and the leading wedge. A big experiment is driven by a single question or hypothesis test, but requires a large scale community investment. [...] A map is a data resource — comprehensive, complete, closed ended — to be used by multiple groups, over a long time, for multiple purposes. The decision to build a map is a cost/benefit calculation, weighed against individual labs who are already making piecemeal maps in an ill coordinated fashion, especially when small groups lack technical expertise to make the map well. A failure mode with a map is to miscalculate the cost/benefit analysis and make a map that too few individual labs will use.

ENCODE and some of its critics have fallen into similar traps. In trying to make the result sound important, ENCODE's publicity spun it retrospectively as a hypothesis test, but ENCODE was not designed to test anything. ENCODE is a map: it should have been published and defended as such. And while its critics argue over an interpretation that wasn't in ENCODE's mission to begin with, ENCODE's planners should also recognize that as ENCODE now moves into a new funding phase, it may be headed for a failure mode in its actual mission. The cost/benefit calculation is rapidly changing. ENCODE's technologies (all based on high throughput sequencing) are now widely and inexpensively available in individual labs.



Bruce Alberts is Editor-in-Chief of *Science*.

EDITORIAL

The End of “Small Science”?

I AM PROMPTED TO WRITE THIS EDITORIAL BY THE RELEASE OF 30 PAPERS THIS MONTH FROM THE ENCODE Project Consortium. This decade-long project involved an international team of 442 scientists who have compiled what is being called an “encyclopedia of DNA elements,” a comprehensive list of functional elements in the human genome. The detailed overview is expected to spur further research on the fundamentals of life, health, and disease. ENCODE exemplifies a “big-science” style of research that continues to sweep the headlines, and the increased efficiency of data production by such projects is impressive. Does this mean that the highly successful “small-science” era of biological research will soon be over? Will government funding increasingly favor big-science projects? I certainly hope that the answer is no.

...

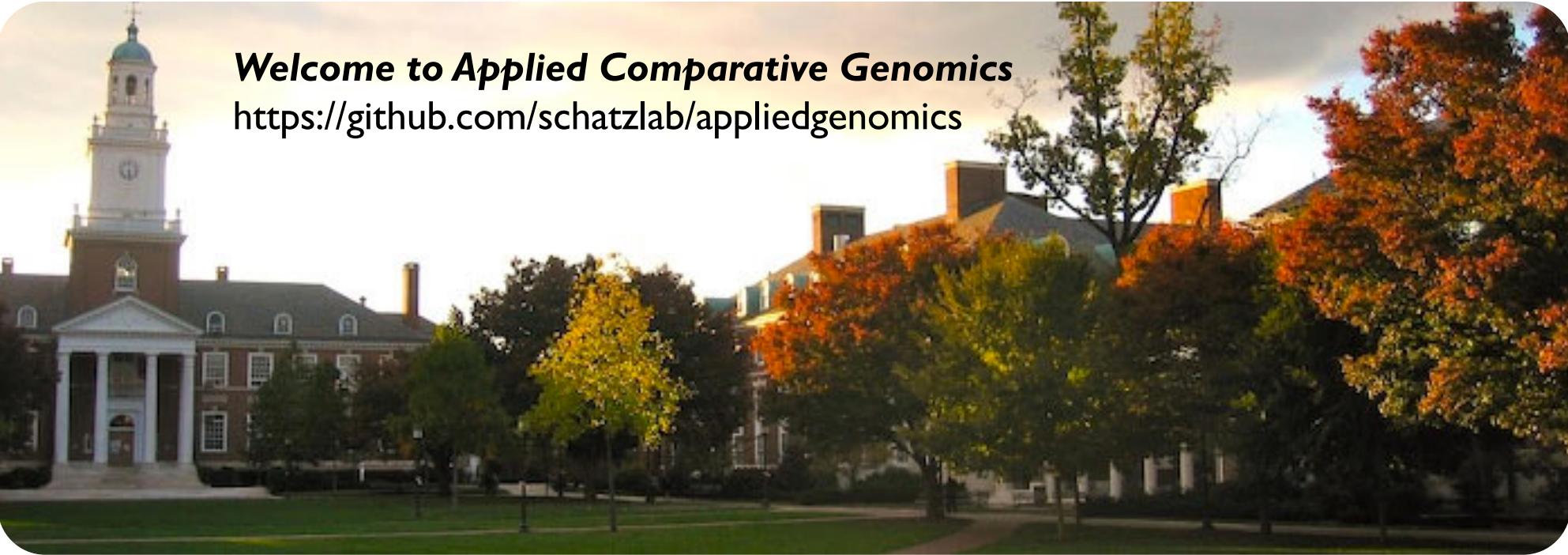
Each year, the amount of factual information that scientists acquire about cells increases and, stimulated by -omics projects, the compilations of data expand at a tremendous rate. But the grand challenges that remain in attaining a deep understanding of the chemistry of life will require going beyond detailed catalogs. Ensuring a successful future for the biological sciences will require restraint in the growth of large centers and -omics-like projects, so as to provide more financial support for the critical work of innovative small laboratories striving to understand the wonderful complexity of living systems.

– Bruce Alberts

10.1126/science.1230529

Next Steps

1. Questions on project?
2. Check out the course webpage



Welcome to Applied Comparative Genomics
<https://github.com/schatzlab/appliedgenomics>

Questions?