Jon K. Lærdahl,
Structural Bioinformatics

# Genome browsers

Oslo
universitetssykehus

UiO : **Department of Informatics**
University of Oslo

---

Jon K. Lærdahl,
Structural Bioinformatics
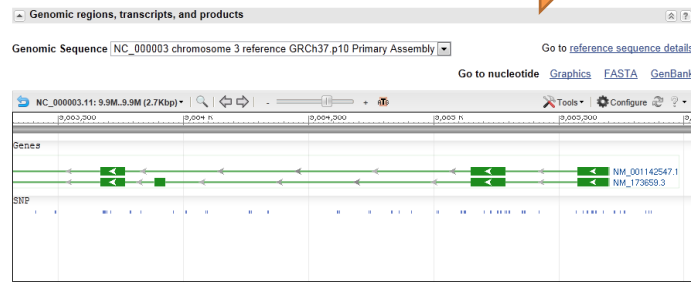
# Genome browsers

- Graphical interface for genomic data
- Shows information from biological databases mapped onto genomic sequence

Genomic coordinates

Genomic regions, transcripts, and products

Genomic Sequence NC_000003 chromosome 3 reference GRCh37.p10 Primary Assembly

Go to reference sequence details

Go to nucleotide   Graphics   FASTA   GenBank

NC_000003.11: 9.9M..9.9M (2.7Kbp)   Tools   Configure

Genes

NM_001142547.1
NM_173659.3

SNP

Various annotations = "tracks"

NCBI Gene database

Oslo
universitetssykehus

UiO : **Department of Informatics**
University of Oslo

---

Jon K. Lærdahl,
Structural Bioinformatics

# UCSC Genome Browser

- Developed and maintained at the University of California, Santa Cruz (UCSC)
- Interactive website
- Access to genome sequence data from
  - Human genome
    - Latest assembly (GRCh37), but also earlier versions
  - Mouse, rat, and approx. 40 other mammals
  - Chicken, turkey, reptiles, frogs, and fish
  - Insects, nematodes, *S. cerevisiae* and more

Oslo universitetssykehus

UiO : Department of Informatics
University of Oslo

---

Jon K. Lærdahl,
Structural Bioinformatics

# UCSC Genome Browser

## The UCSC Genome Browser database: extensions and updates 2011

Timothy R. Dreszer[1,*], Donna Karolchik[1,*], Ann S. Zweig[1], Angie S. Hinrichs[1], Brian J. Raney[1], Robert M. Kuhn[1], Laurence R. Meyer[1], Mathew Wong[1], Cricket A. Sloan[1], Kate R. Rosenbloom[1], Greg Roe[1], Brooke Rhead[1], Andy Pohl[1,2], Venkat S. Malladi[1], Chin H. Li[1], Katrina Learned[1], Vanessa Kirkup[1], Fan Hsu[1], Rachel A. Harte[1], Luvina Guruvadoo[1], Mary Goldman[1], Belinda M. Giardine[3], Pauline A. Fujita[1], Mark Diekhans[1], Melissa S. Cline[1], Hiram Clawson[1], Galt P. Barber[1], David Haussler[1,4] and W. James Kent[1]

[1]Center for Biomolecular Science and Engineering, School of Engineering, University of California Santa Cruz (UCSC), Santa Cruz, CA 95064, USA, [2]Centre for Genomic Regulation (CRG), Barcelona, Spain, [3]Center for Comparative Genomics and Bioinformatics, Huck Institutes of the Life Sciences, Pennsylvania State University, University Park, PA 16802 and [4]Howard Hughes Medical Institute, UCSC, Santa Cruz, CA 95064, USA

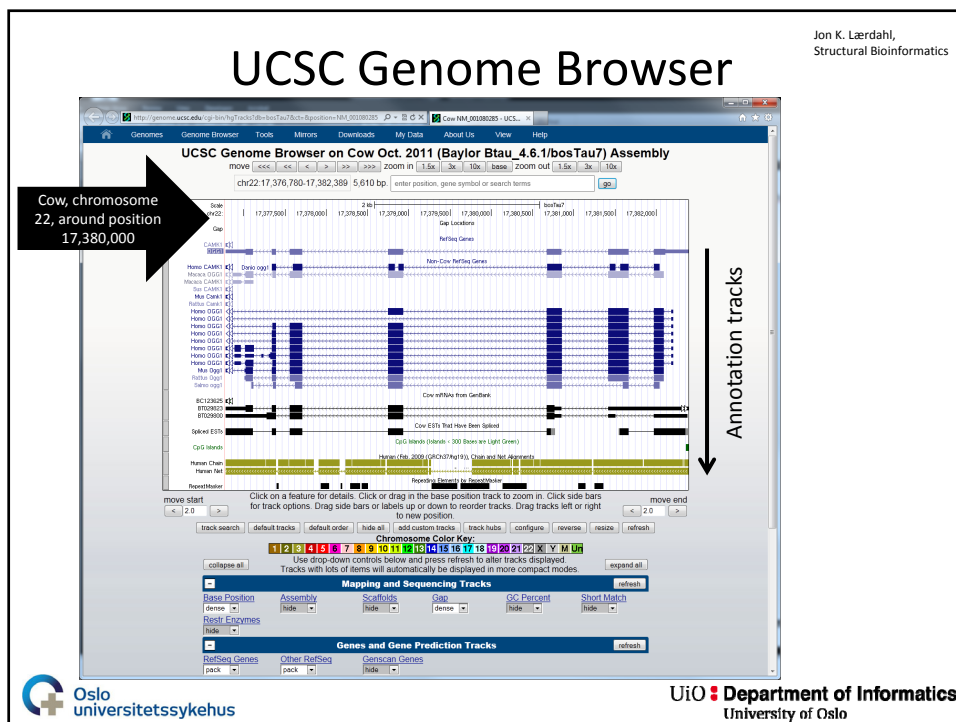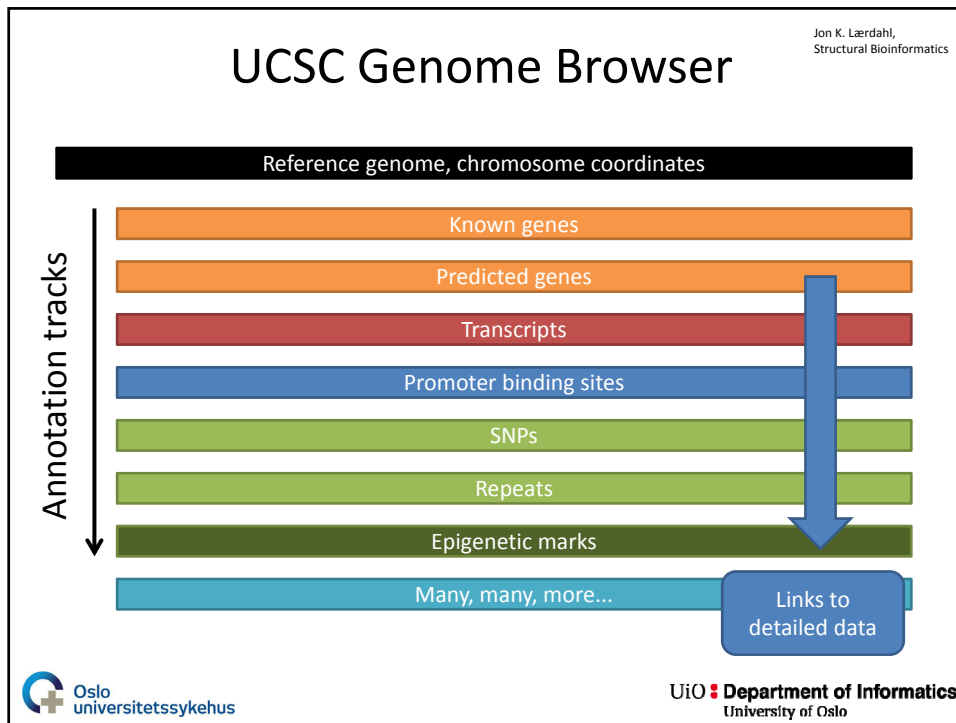Received September 15, 2011; Revised October 18, 2011; Accepted October 25, 2011

http://genome.ucsc.edu

**ABSTRACT**

The University of California Santa Cruz Genome Browser (http://genome.ucsc.edu) offers online public access to a growing database of genomic sequence and annotations for a wide variety of organisms. The Browser is an integrated tool set for visualizing, comparing, analyzing and sharing both publicly available and user-generated genomic data sets. In the past year, the local database has been updated with four new species assemblies, and we anticipate another four will be released by the end

**INTRODUCTION**

The University of California Santa Cruz (UCSC) Genome Browser (1,2) at http://genome.ucsc.edu is a web-based set of tools providing access to a database of genome sequence and annotations for visualization, comparison and analysis by the scientific, medical and academic communities. Our primary mission is to provide timely and convenient open access to high-quality human genome sequence and annotations in a framework that enables easy exploration from genome-wide down to the base level. Annotation datasets, or 'tracks', on the human genome cover conservation and evolutionary compari-

Oslo universitetssykehus

UiO : Department of Informatics
University of Oslo

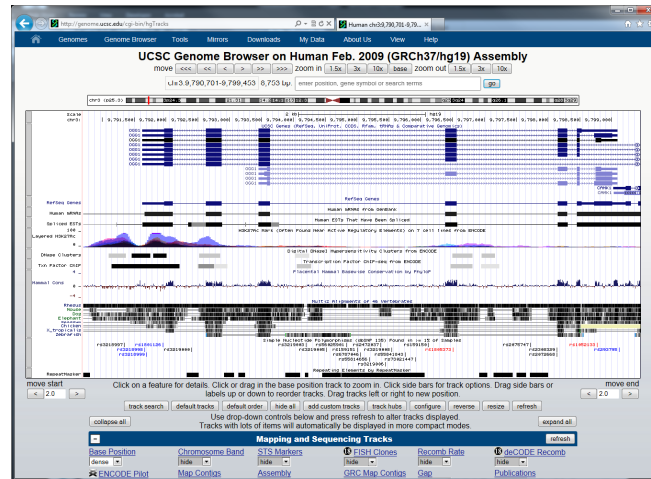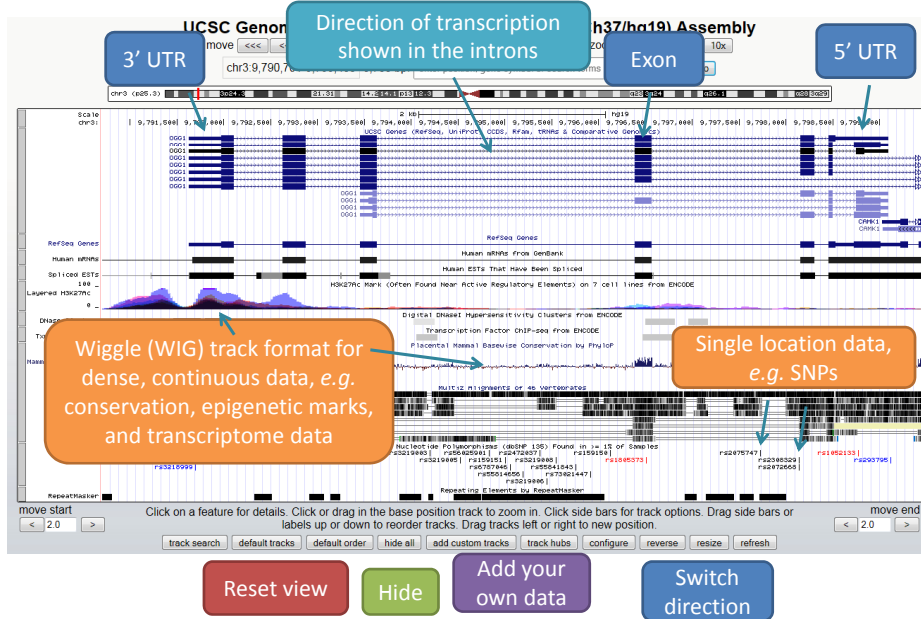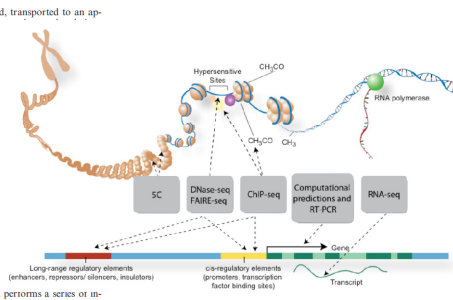# UCSC Genome Browser

Jon K. Lærdahl,
Structural Bioinformatics

Access to the databases and tools

Start here

http://genome.ucsc.edu

General information

News, updates, announcements

Oslo universitetssykehus

UiO : Department of Informatics
University of Oslo

---

# UCSC Genome Browser

Jon K. Lærdahl,
Structural Bioinformatics

Examples of searching options – correct query format

Oslo universitetssykehus

UiO : Department of Informatics
University of Oslo

# UCSC Genome Browser
# brief demo

Jon K. Lærdahl,
Structural Bioinformatics



Oslo universitetssykehus

UiO : **Department of Informatics**
University of Oslo

# Different kinds of data

Jon K. Lærdahl,
Structural Bioinformatics



3' UTR

Direction of transcription shown in the introns

Exon

5' UTR

Wiggle (WIG) track format for dense, continuous data, *e.g.* conservation, epigenetic marks, and transcriptome data

Single location data, *e.g.* SNPs

Reset view

Hide

Add your own data

Switch direction

## Slide: ENCODE data in UCSC

Jon K. Lærdahl,
Structural Bioinformatics

# ENCODE data in UCSC

**ENCODE whole-genome data in the UCSC genome browser (2011 update)**

Brian J. Raney[1,*], Melissa S. Cline[1], Kate R. Rosenbloom[1], Timothy R. Dreszer[1], Katrina Learned[1], Galt P. Barber[1], Laurence R. Meyer[1], Cricket A. Sloan[1], Venkat S. Malladi[1], Krishna M. Roskin[1], Bernard B. Suh[1], Angie S. Hinrichs[1], Hiram Clawson[1], Ann S. Zweig[1], Vanessa Kirkup[1], Pauline A. Fujita[1], Brooke Rhead[1], Kayla E. Smith[1], Andy Pohl[1], Robert M. Kuhn[1], Donna Karolchik[1], David Haussler[1,2] and W. James Kent[1]

[1]Center for Biomolecular Science and Engineering, School of Engineering and [2]Howard Hughes Medical Institute, University of California Santa Cruz (UCSC), Santa Cruz, CA 95064, USA

ABSTRACT

The ENCODE project is an international consortium with a goal of cataloguing all the functional elements in the human genome. The ENCODE Data Coordination Center (DCC) at the University of California, Santa Cruz serves as the central repository for ENCODE data. In this role, the DCC offers a collection of high-throughput, genome-wide data generated with technologies such as ChIP-Seq, RNA-Seq, DNA digestion and others. This data helps illuminate transcription factor-binding sites, histone marks, chromatin accessibility, DNA methylation, RNA expression, RNA binding and other cell-state indicators. It includes sequences with quality scores, alignments, signals calculated from the alignments, and in most cases, element or peak calls calculated from the signal data. Each data set is available for visualization and download via the UCSC Genome Browser (http://genome.ucsc.edu/). ENCODE data can also be retrieved using a metadata system that captures the experimental parameters of each assay. The ENCODE web portal at UCSC (http://encodeproject.org/) provides information about the ENCODE data and links for access.

into RNA, which might be spliced, transported to an appropriate cellular compartmen
proteins. This process is regulated
DNA methylation, chromatin n
transcription factors to the DN
factors to the RNA and RNA tran
itable traits are determined as m
regulation as differences in gene e

The goal of the ENCODE proj
tional elements in the human ge
processes, through direct measur
genomic technologies and detai
ENCODE began with a pilot ph
of the genome (3), and scaled
analysis in 2007.

The role of the ENCODE Da
(DCC) is to organize and display
labs in the consortium, and to en
specific quality standards when it
Before a lab submits any data, th
a data agreement that defines the
and associated metadata. The D
data to ensure consistency with
loads the data onto a test serve
tion, and coordinates with the l
into a consistent set of tracks. W
the DCC Quality Assurance team performs a series of in-

http://genome.ucsc.edu/ENCODE/aboutScaleup.html

Oslo universitetssykehus

UiO : **Department of Informatics**
**University of Oslo**

---

## Slide: Ensembl Genome Browser

Jon K. Lærdahl,
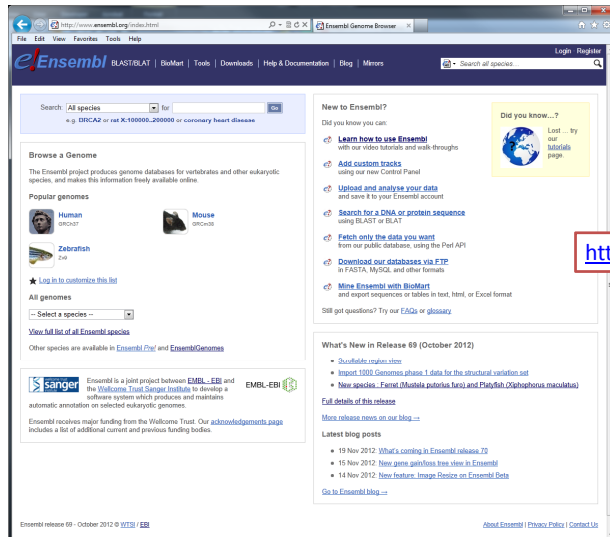Structural Bioinformatics

# Ensembl Genome Browser

- Joint project between EMBL-EBI and the Wellcome Trust Sanger Institute
- Central resource for studying genomes of vertebrates
  - Mainly chordates, but some few extra (*e.g. C. elegans and S. cerevisiae*)
  - Updated several times a year with new genome assemblies and new species
  - Annotations of genomes (*e.g.* genes and their splice variant, SNPs) added by the Ensembl pipeline
  - Automatic gene prediction (with or without experimental evidence) & some curator input

Oslo universitetssykehus

UiO : **Department of Informatics**
**University of Oslo**

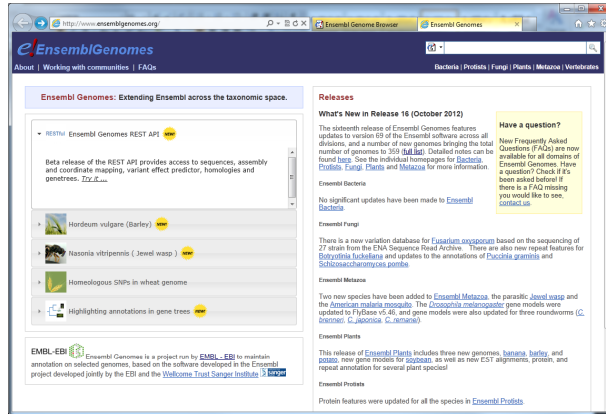Jon K. Lærdahl,
Structural Bioinformatics

# *EnsemblGenomes*



• Bacteria, protists, fungi, plants and other metazoa (359 genomes)

Oslo
universitetssykehus

UiO **: Department of Informatics**
**University of Oslo**

---

Jon K. Lærdahl,
Structural Bioinformatics

# Ensembl Genome Browser



http://www.ensembl.org

# Explore!

Oslo
universitetssykehus

UiO **: Department of Informatics**
**University of Oslo**

Jon K. Lærdahl,
Structural Bioinformatics

# Now something different!

Not a genome browser!

Oslo universitetssykehus

UiO **: Department of Informatics**
University of Oslo

---



Jon K. Lærdahl,
Structural Bioinformatics

- Galaxy is a platform (open, web-based) for computational medical projects and bioinformatics
  - Accessible: Not necessary to know programming, Unix, or how to install programs
  - Reproducible: You can build and store complete workflows, pipelines, and the full computational analysis
  - Transparent: Users can publish and share whole worksflows
- A bioinformatics workflow management system

Oslo universitetssykehus

UiO **: Department of Informatics**
University of Oslo

# Reproducibility

Jon K. Lærdahl,
Structural Bioinformatics

## Repeatability of published microarray gene expression analyses

John P A Ioannidis[1–3], David B Allison[4], Catherine A Ball[5], Issa Coulibaly[4], Xiangqin Cui[4], Aedín C Culhane[6,7], Mario Falchi[8,9], Cesare Furlanello[10], Laurence Game[11], Giuseppe Jurman[10], Jon Mangion[11], Tapan Mehta[4], Michael Nitzberg[5], Grier P Page[4,12], Enrico Petretto[11,13] & Vera van Noort[14]

Given the complexity of microarray-based gene expression studies, guidelines encourage transparent design and public data availability. Several journals require public data deposition and several public databases exist. However, not all data are publicly available, and even when available, it is unknown whether the published results are reproducible by independent scientists. Here we evaluated the replication of data analyses in 18 articles on microarray-based gene expression profiling published in *Nature Genetics* in 2005–2006. One table or figure from each article was independently evaluated by two teams of analysts. We reproduced two analyses in principle and six partially or with some discrepancies; ten could not be reproduced. The main reason for failure to reproduce was data unavailability, and discrepancies were mostly due to incomplete data annotation or specification of data processing and analysis. Repeatability of published microarray studies is apparently limited. More strict publication rules enforcing public data availability and explicit description of data processing and analysis should be considered.

**Figure 1** Summary of the efforts to replicate the published analyses.

Oslo universitetssykehus

UiO **Department of Informatics** University of Oslo

---

# Reproducibility

Jon K. Lærdahl,
Structural Bioinformatics

Raw data

Was it version 2.4 or 3.1?

Run program 1

Run program 2

I had to add my own module, but I am not sure I can remember which parts I added

I have it stored somewhere?

Run my own python script

Run program 3

Publication?

This postdoc left the lab. I don't know what happened to that script

Process with perl script I got from someone

Oslo universitetssykehus

UiO **Department of Informatics** University of Oslo

## Slide 1

Jon K. Lærdahl,
Structural Bioinformatics



Li, M. *et al.* Science **33**, 53 (2011), July 2011:
- >10,000 sites in the human genome where an RNA sequence did not match the sequence of the DNA
- Evidence for a new mechanism of RNA editing?
- Central dogma is wrong in 10,000 places?

Oslo
universitetssykehus

UiO **: Department of Informatics**
University of Oslo

## Slide 2

Jon K. Lærdahl,
Structural Bioinformatics

# Galaxy

- Galaxy is a platform (open, web-based) for computational medical projects and bioinformatics
  - Accessible: Not necessary to know programming, Unix, or how to install programs
  - ***Reproducible: You can build and store complete workflows, pipelines, and the full computational analysis***
  - ***Transparent: Users can publish and share whole worksflows***
- By the way, Galaxy is written in Python...
- Developed by the labs of Anton Nekrutenko (Penn State University) and James Taylor (Emory University)

Oslo
universitetssykehus

UiO **: Department of Informatics**
University of Oslo

# Galaxy

Jon K. Lærdahl,
Structural Bioinformatics

- Bioportal -> Galaxy
- Galaxy at UiO submitting jobs to the Abel cluster

Oslo
universitetssykehus

UiO : **Department of Informatics**
University of Oslo

---

# Galaxy

Jon K. Lærdahl,
Structural Bioinformatics

- Can be run on a free public server at Penn State
- You can install Galaxy on your own server or computer cluster (soon on Abel)
- You can run Galaxy in the cloud

Oslo
universitetssykehus

UiO : **Department of Informatics**
University of Oslo

Free public server:

Jon K. Lærdahl,
Structural Bioinformatics



Jon K. Lærdahl,
Structural Bioinformatics

After a while...

Jon K. Lærdahl,
Structural Bioinformatics



Jon K. Lærdahl,
Structural Bioinformatics

Jon K. Lærdahl,
Structural Bioinformatics



Jon K. Lærdahl,
Structural Bioinformatics

Jon K. Lærdahl,
Structural Bioinformatics



- Create a "workflow" from your history
- Edit workflow
    - Edit settings for all "modules"
    - Add new modules
- Share workflow
    - Share on a website?
    - As supplementary material for publication?

# Galaxy

Jon K. Lærdahl,
Structural Bioinformatics

- A *very* simple demonstration
- Galaxy can quite easily answer questions like:
  - Which exon in the human genome contains the largest number of SNPs?
  - How many exons are there on mouse chromosome 1?
- Very good for making pipeline to analyze HTS data
  - ChIP-seq, RNA-seq etc
- If you are doing this kind of work, Galaxy might be something to consider!
- Try it yourself? Check out the Galaxy 101 screencast?

Oslo
universitetssykehus

UiO **: Department of Informatics**
University of Oslo