**When do Scientists Change their Minds?**
*Week 4 – The Human Genome*

EGMT-1520  Mon, Feb 7, 2022
Bill Pearson  wrp@virginia.edu

Overview of this session:
- What is a genome?
- The human genome project
  – the beginnings (mapping, cloning)
  – the end (industrial sequencing, shotgun genomes)
  – the sequel (Next Generation sequencing)
- The human genome discoveries
  – number of genes
  – gene organization
  – genome conservation
- Browsing the genome

1

# For Wednesday:

Human genome lab (in groups) –
1. look up a gene in the human genome
2. characterize the gene
   a. identify beginning, end
   b. count the number of exons
   c. count the number of mRNA isoforms
3. find the nearest gene "upstream" and "downstream"
4. characterize the "upstream" or "downstream" region
   a. how conserved is the upstream/downstream region compared to the exons in your gene from humans to chimps (5 Mya)?
   b. from humans to mouse (80 Mya)?  Is the conservation uniform?
   c. what features are annotated in this region?  repeated sequences? other conserved regions?

2

1

## For Monday:

Repeat the human genome lab on a different gene.

1. report the name of the gene, and its chromosome location.  Submit the URL of the UCSC genome browser page that shows the gene.
2. characterize the gene
    a. report the the length of the gene
    b. Is the gene on the forward or reverse strand?
    c. report the number of exons
3. report the name and coordinates of the  nearest gene "upstream" and "downstream"
    a. Determine whether the gene is on the same strand, (forward/reverse) or on the opposite strand.
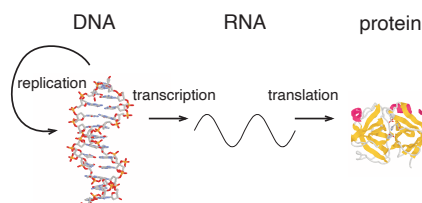
fasta.bioch.virginia.edu/egmt1520                      3

3

## What does a gene look like?

- Central dogma:

DNA          RNA          protein

replication        transcription        translation



- Parts of a gene:
    – start of transcription  (beginning of the mRNA)
    – start of translation  (beginning of the protein)
    – end of translation  (end of the protein)
    – poly-A addition site (end of processed mRNA)
    – end of transcription
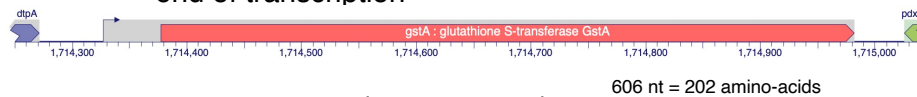
fasta.bioch.virginia.edu/egmt1520                      4

4

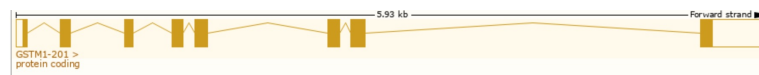## Prokaryotes (bacteria, archea) vs Eukaryotes

- Parts of a gene (prokaryotes, eukaryotes):
  – start of transcription  (beginning of the mRNA)
  – start of translation  (beginning of the protein)
  – end of translation  (end of the protein)
  – poly-A addition site (end of processed mRNA)
  – end of transcription



606 nt = 202 amino-acids

- Parts of a gene (eukaryotes):
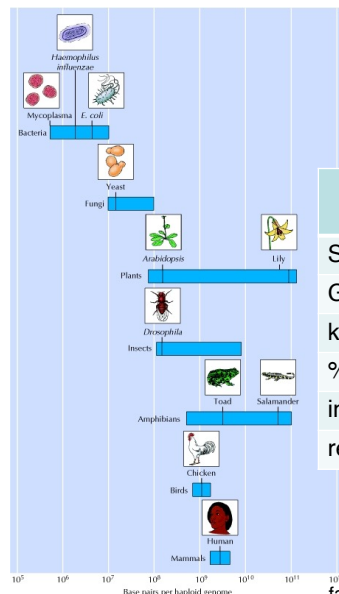  – introns and exons



gene: 5950 nt, mRNA: 1164 nt, protein: 218 aa

fasta.bioch.virginia.edu/egmt1520

5

5

## What is in a genome?



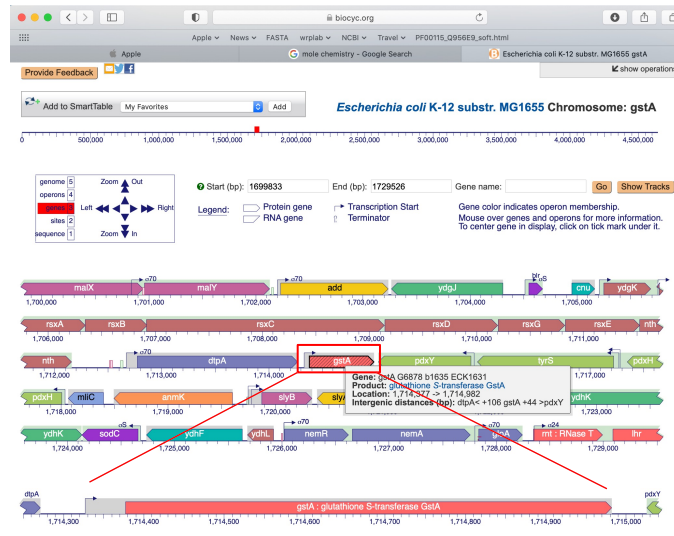|  | E. col | Plas. | Yeast | Plant (ARATH) | Homo |
|---|---|---|---|---|---|
| Size(Mb) | 4.64 | 22.8 | 12.5 | 115 | 3289 |
| Genes | 4288 | 5268 | 5770 | 25.5K | ~25K |
| kb/Gene | 0.95 | 4.34 | 2.09 | 4.53 | 27 |
| %coding | 87.8 | 52.6 | 70.5 | 28.8 | 1.3 |
| introns | 0 | 7406 | 272 | 107K | 53K |
| repeat% | <1 | <1 | 2.4 | 15 | 46 |

Pevsner, Table 16-1

Cooper, GM (2000) The Cell: A Molecular
Approach. 2nd edition. Fig 4.1

fasta.bioch.virginia.edu/egmt1520

6

6

Gene and genome complexity – E. coli
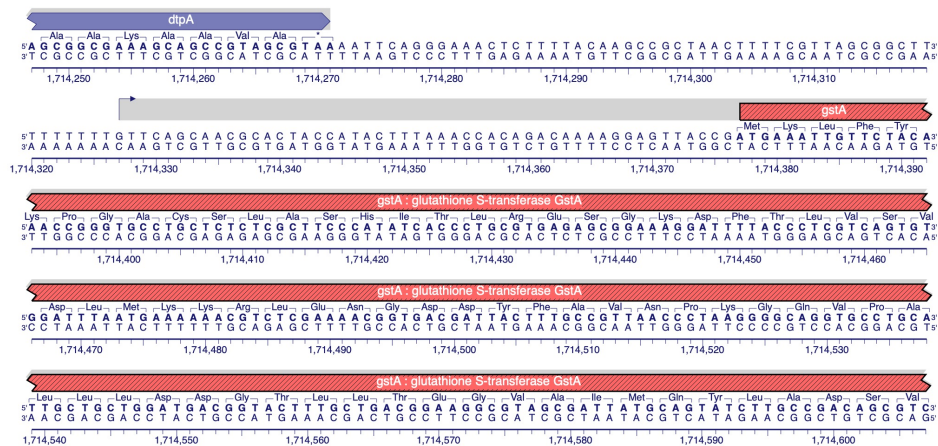(4288 genes, 4.64 Mbase, 87.8% protein coding)
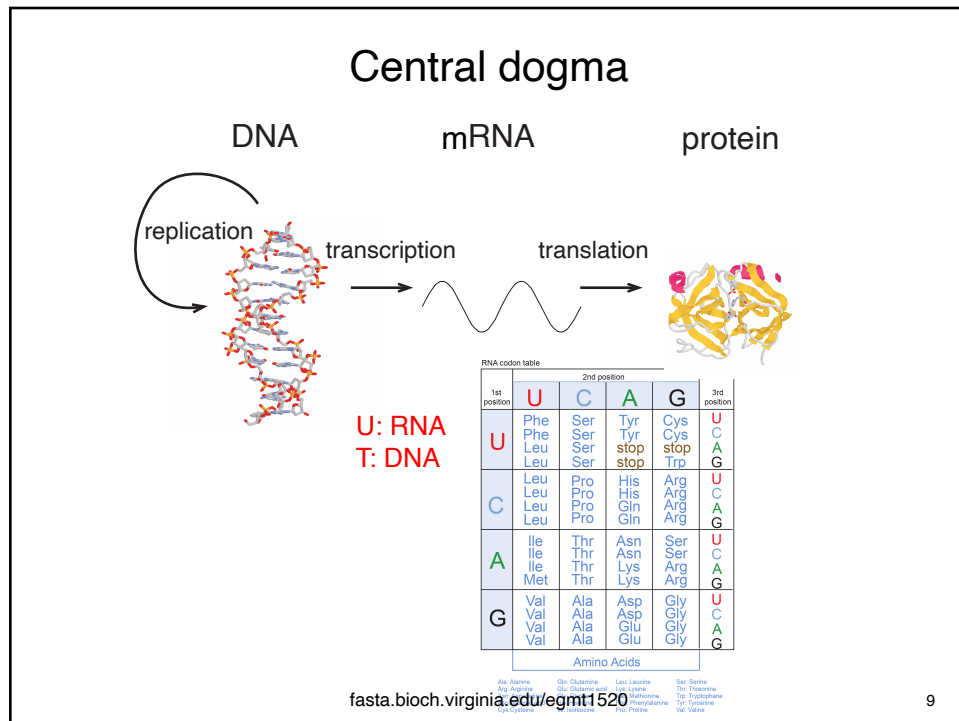
fasta.bioch.virginia.edu/egmt1520

biocyc.org

7

7



Gene and genome complexity – E. coli
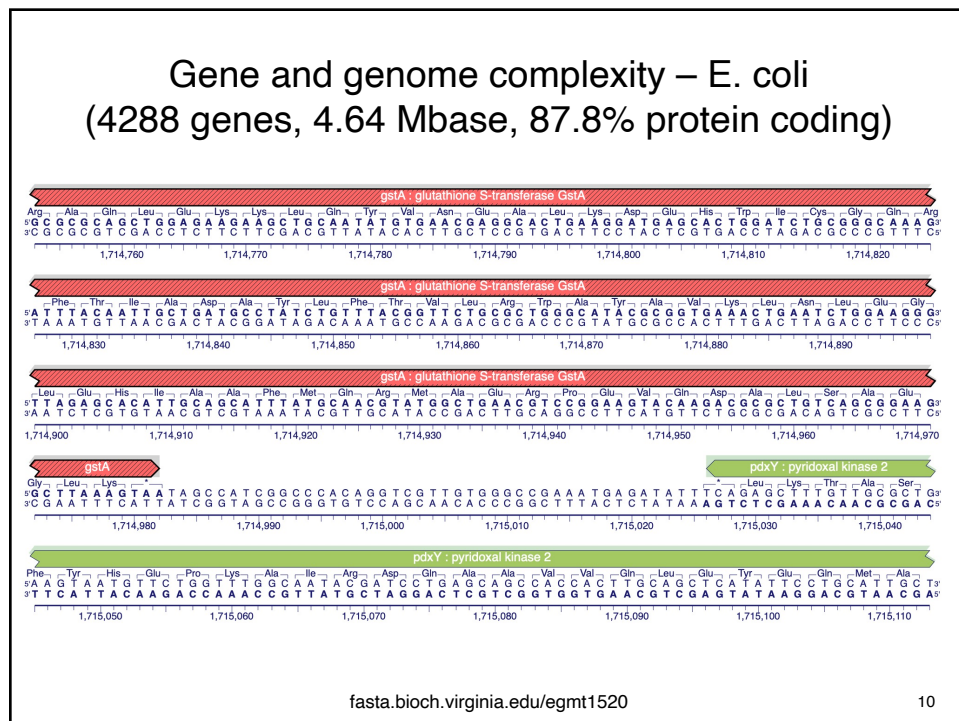(4288 genes, 4.64 Mbase, 87.8% protein coding)

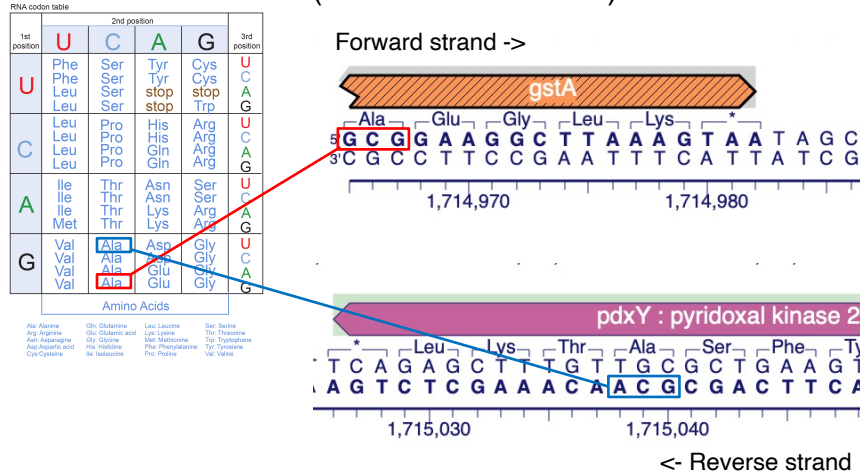fasta.bioch.virginia.edu/egmt1520

8

8

Central dogma

DNA          mRNA          protein

replication

transcription          translation

U: RNA
T: DNA

fasta.bioch.virginia.edu/egmt1520

9

Gene and genome complexity – E. coli
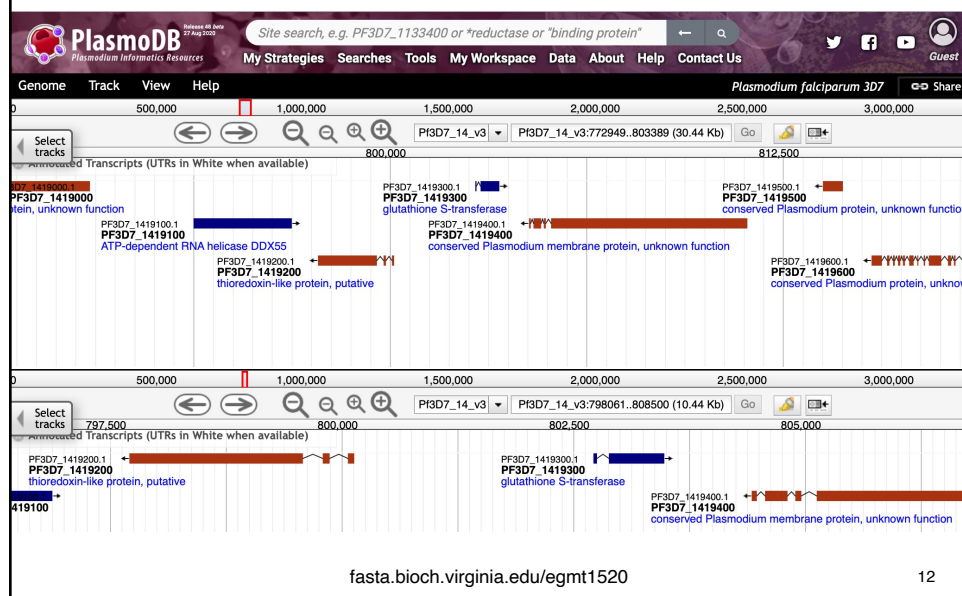(4288 genes, 4.64 Mbase, 87.8% protein coding)
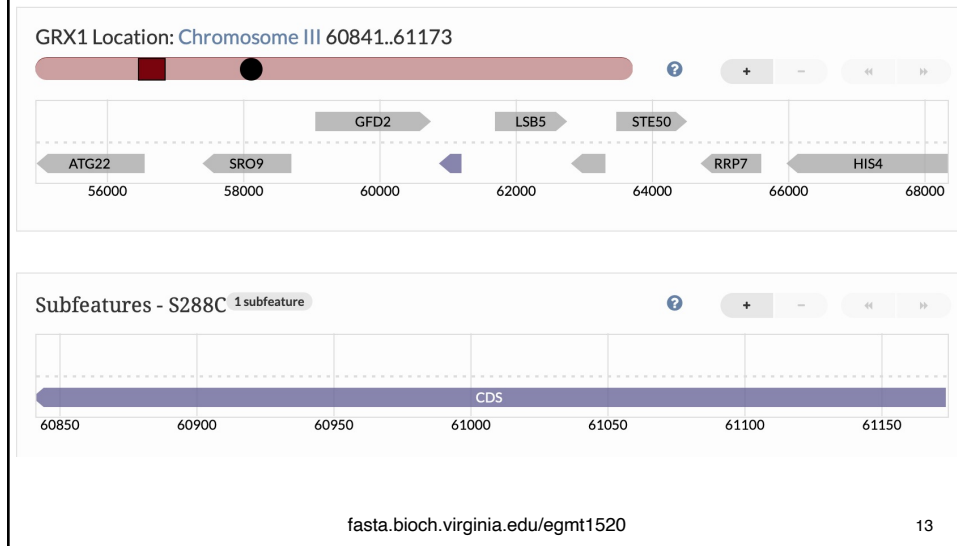
fasta.bioch.virginia.edu/egmt1520
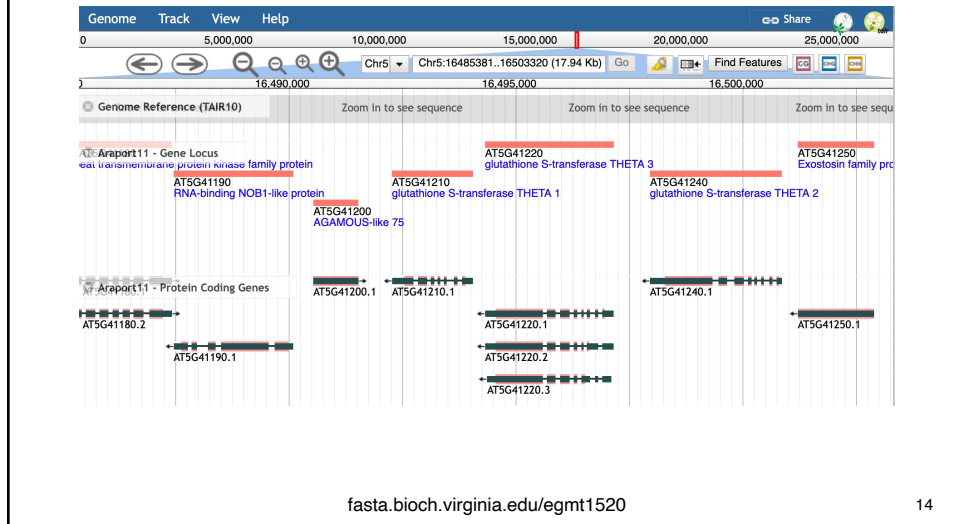
10

11



12

Gene and genome complexity – Yeast
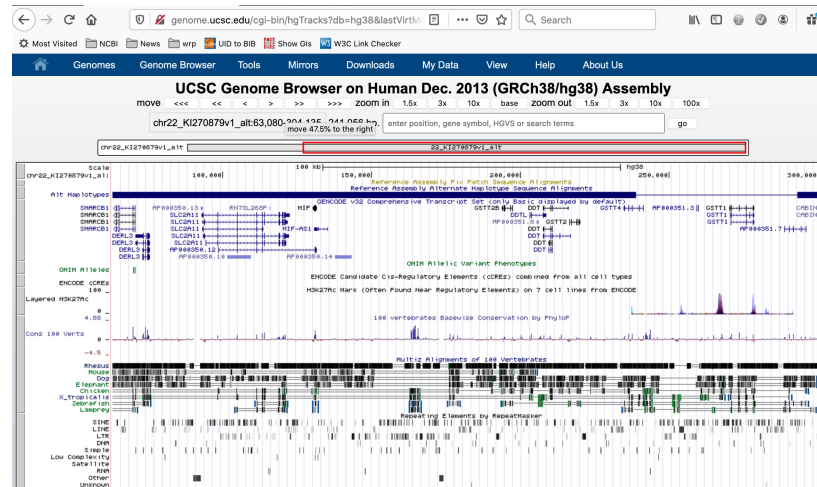(5770 genes, 12.5 Mbase, 70.5% protein coding)

13



Gene and genome complexity –
Arabidopsis (mouse ear cress)
(22,500 genes, 115 Mbase, 28 % protein coding)

14

## Gene and genome complexity – Human
## (~20,000 genes, 3,000 Mbase, 1.5% protein coding)



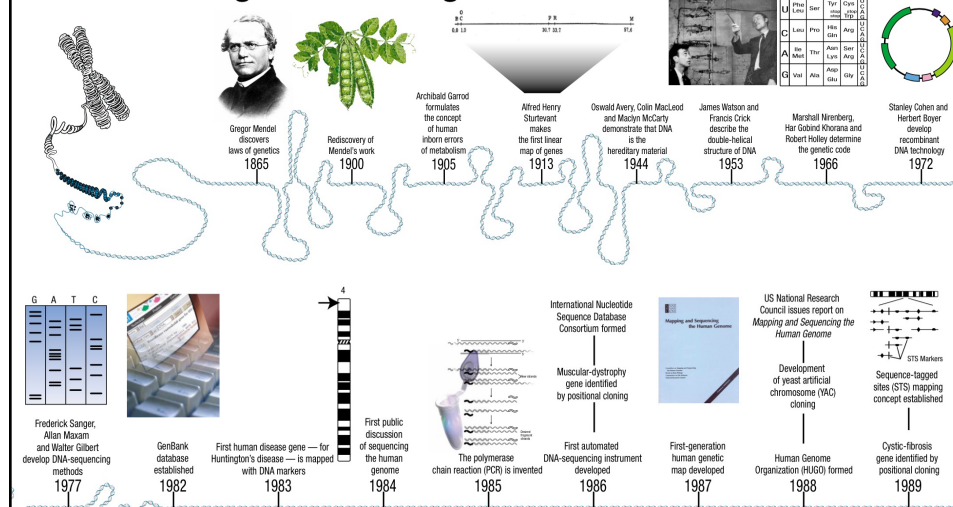fasta.bioch.virginia.edu/egmt1520

15

15

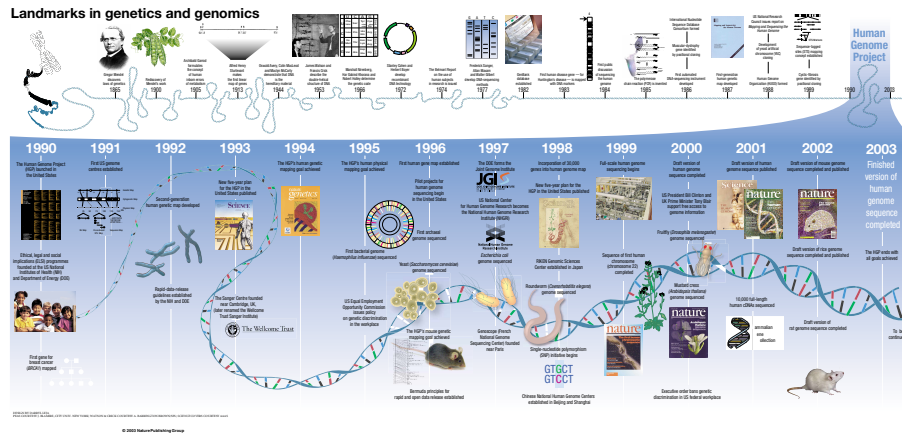## Landmarks in genetics and genomics



Collins, Nature (2003) 422:835

fasta.bioch.virginia.edu/egmt1520
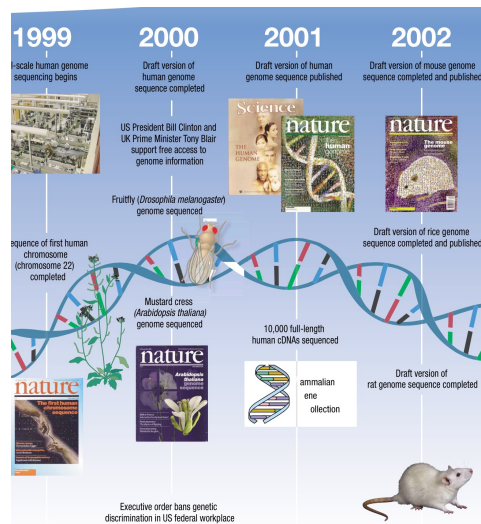
16

16

# A history of genomes



Collins, Nature (2003) 422:835

fasta.bioch.virginia.edu/egmt1520

17

17

# The Human Genome Project



Collins, Nature (2003) 422:835

fasta.bioch.virginia.edu/egmt1520

18

18

Sequencing capacity (2004)
40 million lanes/year ~ 8 billion bases – Whitehead Inst. (1 bacterial genome/day)
>40 billion bases/year, world-wide
(1000 bacterial genomes/year; 1 mammalian genome/year)

Sequencing capacity (2011) – Illumina sequencing 200 billion bases/week/machine, ~30,000 human genomes/year
Sequencing capacity (2015) – at least 300,000 human genomes/year

fasta.bioch.virginia.edu/egmt1520          19

19

# The human genome – initial insights

1. There were reported to be about 30,000 to 40,000 predicted protein-coding genes in the human genome. Currently, ENSEMBL reports 20,300 protein coding genes. Similar to Arabidopsis (plant, 26,000 genes) and pufferfish (21,000 genes), and marginally more genes than are found in many nematode and insect genomes (14,000).

2. More than 98% of the human genome does not code for genes. Much of this genomic landscape is occupied by repetitive DNA elements such as long interspersed elements (LINEs) (20%), short interspersed elements (SINEs) (13%), long terminal repeat (LTR) retrotransposons (8%), and DNA transposons (3%). Thus half the human genome is derived from transposable elements.

3. The mutation rate is about twice as high in male meiosis than in female meiosis. This suggests that most mutation occurs in males.

4. More than 1.4 million single nucleotide polymorphisms (SNPs) were identified. SNPs are single nucleotide variations that occur once every 100 to 300 base pairs (bp). 36 million in Oct., 2014

5. Comparative sequencing of the mouse genome suggests that only about 5% of the human genome is under evolutionary selection.

fasta.bioch.virginia.edu/egmt1520          20

20

## The human genome
## (GRCh38.p13 Dec. 2013)

**Human assembly and gene annotation**

www.ensembl.org

**Assembly**

This site provides a data set based on the December 2013 *Homo sapiens* high coverage assembly GRCh38 from the Genome Reference Consortium . This assembly is used by UCSC to create their hg38 database. The data set consists of gene models built from the genewise alignments of the human proteome as well as from alignments of human cDNAs using the cDNA2genome model of exonerate.

This release of the assembly has the following properties:

- contig length total 3.4 Gb.
- chromosome length total 3.1 Gb (excluding haplotypes).

It also includes 261 alt loci scaffolds, mainly in the LRC/KIR complex on chromosome 19 (35 alternate sequence representations) and the MHC region on chromosome 6 (7 alternate sequence representations).

Watch a video on YouTube about patches and haplotypes in the Human genome.

**Patches**

As the GRC maintains and improves the assembly, patches are being introduced. Currently, assembly patches are of two types:

- Novel patch: new sequences that add alternative sequence at a loci and will remain as haplotypes in the next major assembly release by GRC
- Fix patch: sequences that correct the reference sequence and will replace the given region of the reference assembly at the next major assembly release by GRC.

**Other assemblies**

GRCh37 Full Feb 2014 archive with BLAST, VEP and BioMart ⌄  Go

**Statistics**

**Summary**

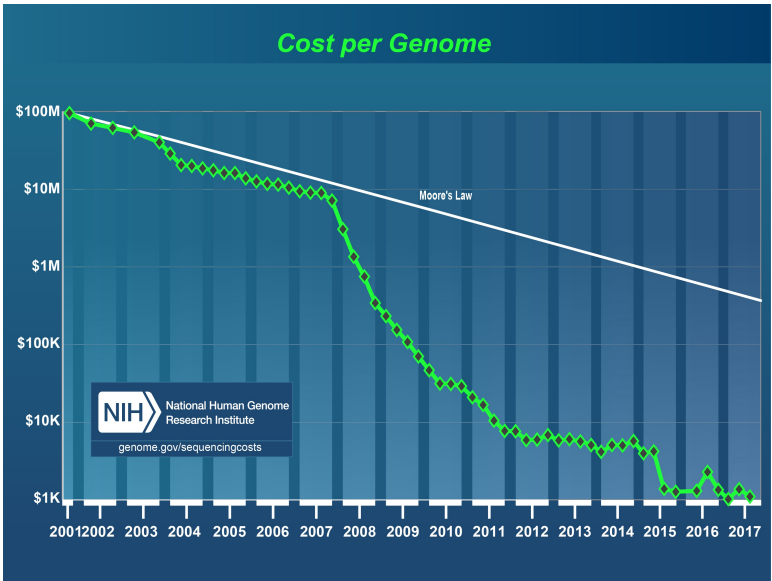| | |
|---|---|
| Assembly | GRCh38.p13 (Genome Reference Consortium Human Build 38), INSDC Assembly GCA_000001405.28 , Dec 2013 |
| Base Pairs | 4,537,931,177 |
| Golden Path Length | 3,096,649,726 |
| Annotation method | Full genebuild |
| Genebuild started | Jan 2014 |
| Genebuild released | Jul 2014 |
| Genebuild last updated/patched | Jun 2019 |
| Database version | 101.38 |
| Gencode version | GENCODE 34 |

**Gene counts (Primary assembly)**

| | |
|---|---|
| Coding genes | 20,440 (incl 633 readthrough) |
| Non coding genes | 23,995 |
| Small non coding genes | 4,867 |
| Long non coding genes | 16,907 (incl 306 readthrough) |
| Misc non coding genes | 2,221 |
| Pseudogenes | 15,222 (incl 6 readthrough) |
| Gene transcripts | 229,649 |

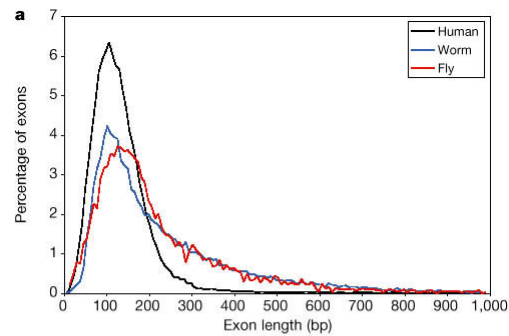21

## Sequencing Costs



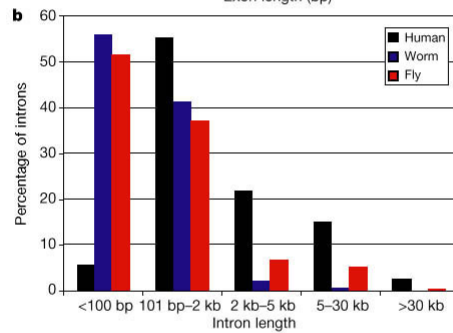fasta.bioch.virginia.edu/egmt1520

22

22

## Exon Length

## Intron Length
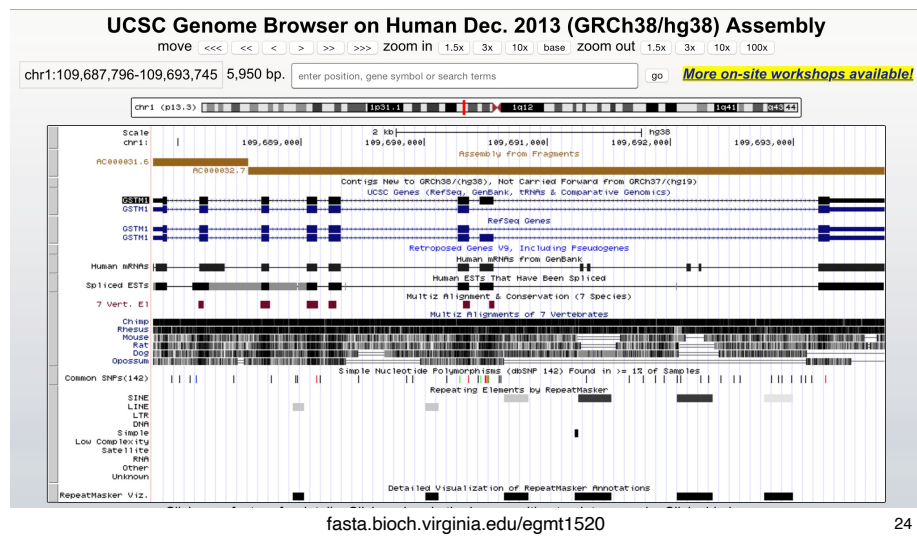
Lander *Nature* **409,** 860–921 (2001), Fig. 35

fasta.bioch.virginia.edu/egmt1520

23

23

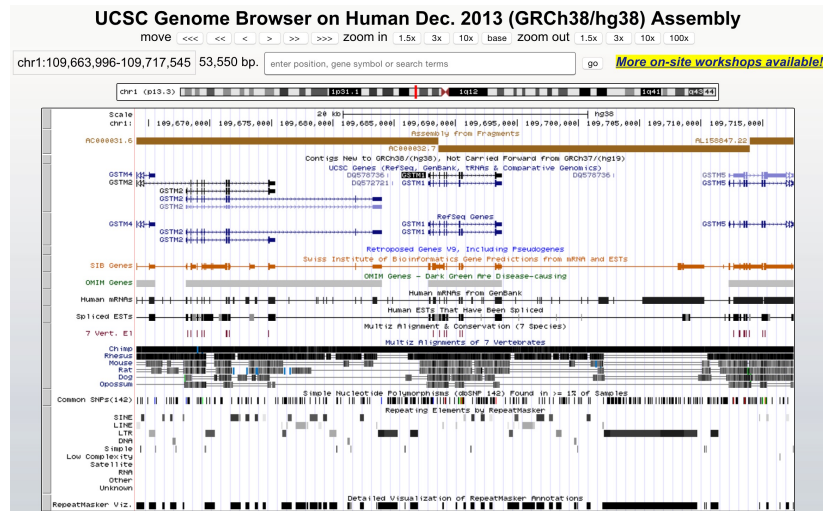# Retrieving the data: Genome Browsers (UCSC)

fasta.bioch.virginia.edu/egmt1520

24

24

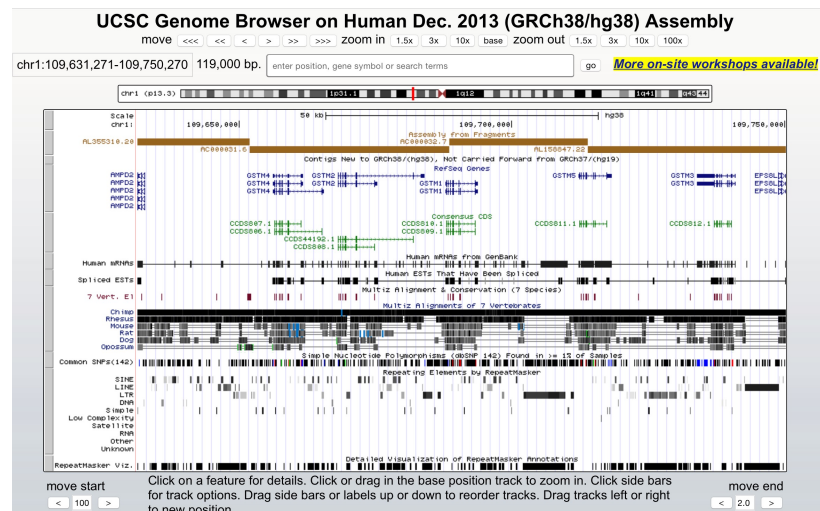## Genome Browsers (genome.ucsc.edu)



fasta.bioch.virginia.edu/egmt1520    25

25

## Genome Browsers (UCSC)



fasta.bioch.virginia.edu/egmt1520    26

26

## The human genome – initial insights

1. There were reported to be about 30,000 to 40,000 predicted protein-coding genes in the human genome. Currently, ENSEMBL reports 20,300 protein coding genes. Similar to Arabidopsis (plant, 26,000 genes) and pufferfish (21,000 genes), and marginally more genes than are found in many nematode and insect genomes (14,000).
2. More than 98% of the human genome does not code for genes. Much of this genomic landscape is occupied by repetitive DNA elements such as long interspersed elements (LINEs) (20%), short interspersed elements (SINEs) (13%), long terminal repeat (LTR) retrotransposons (8%), and DNA transposons (3%). Thus, half the human genome is derived from transposable elements.
3. The mutation rate is about twice as high in male meiosis than in female meiosis. This suggests that most mutation occurs in males.
4. More than 1.4 million single nucleotide polymorphisms (SNPs) were identified. SNPs are single nucleotide variations that occur once every 100 to 300 base pairs (bp). 36 million in Oct., 2014
5. Comparative sequencing of the mouse genome suggests that only about 5% of the human genome is under evolutionary selection.

fasta.bioch.virginia.edu/egmt1520                27

27

## For Wednesday:

Human genome lab (in groups) –
Do the human genome exploration lab exercises.

1. characterize a gene
   a. identify beginning, end
   b. count the number of exons
   c. count the number of mRNA isoforms
2. find the nearest gene "upstream" and "downstream"
3. characterize the "upstream" or "downstream" region
   a. how conserved is the upstream/downstream region compared to the exons in your gene from humans to chimps (5 Mya)?
   b. from humans to mouse (80 Mya)?  Is the conservation uniform?
   c. what features are annotated in this region?  repeated sequences? other conserved regions?

fasta.bioch.virginia.edu/egmt1520                28

28

# For Monday:

Repeat the human genome lab on a different gene.

1. report the name of the gene, and its chromosome location.  Submit the URL of the UCSC genome browser page that shows the gene.
2. characterize the gene
   a. report the the length of the gene
   b. Is the gene on the forward or reverse strand?
   c. report the number of exons
3. report the name and coordinates of the  nearest gene "upstream" and "downstream"
   a. Determine whether the gene is on the same strand, (forward/reverse) or on the opposite strand.

fasta.bioch.virginia.edu/egmt1520                    29

29