# Genomics I

Biomedical Data Science: Mining and Modeling
CB&B 752 • MB&B 452
Matt Simon
January 15, 2020

# What is genomics?

1. The **global** study of how biological **information** is encoded in
   genome sequence

   Genes
   Regulatory sequences
   Genetic variation

2. How this information is **read out** to produce distinct **biological
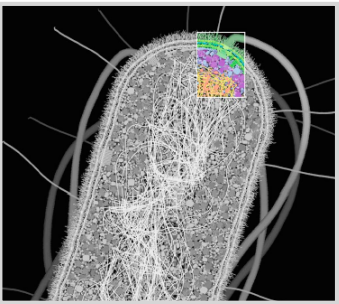   outcomes**

   Gene expression and regulation
   Cellular identity, differentiation and development
   Phenotypic variation among individuals and species

In practice, many experiments that involve
**deep sequencing** are considered genomics.

TCAAATTCGACTGAGAATAAAACAGACACTAA
CTTCTTTCACTTCTTACCTGTCAATGTTATTAA
CATTCAATAAATATTTTTTAGAATAATAAGTCC
ATAAAACAGACACAAACAAGTAAATAAAGTTA
AAATCACATTAATTCCAACATGCAAAGAGGAA
AATGTGGCCAACATGCAAAGAGGAAATCTC
TAGGATAAGGATAATATACAGAGAACATGCCA
TTCACTTCTTACCTGTCAATGTTATTAATATTT
AATAAATATTTTTTAGAATAATAAGTCCCAGGC
ACAGACACTAAACAAGTAAATAAAGTTAATTT
TGGGATTGGAAGACCTCTCTGAGATTAGTGT
GAATGAGCTGGATATACTCAAGGAAGAAGA
TTTTAGGAACAATAAATCACATTAATTCCTTAT
GCACAAGACCAGTATTATGTTCTAGGCATTG
TTCAAGTTGTAATTGATGCTACTATGGAAAA
CACATTAATTCCTTATCTCATGTGAAATTCAT
ATGTTCTAGGCATTGGGGATACCATGTTCAC
GCTATCCCAGGCACAAGACCAGTATTATGTT
CATTCGTTTATCAGAGGCCAAATGTTTTTCTT
AGTTGTATTATTAGAAACTGAGGGCTAAAAC
AAACAAAGACTGTTACTATGGAAAAATGAAA
TAATTCCTTATCTCATGTGAAATTCATATTTA
CTAGGCATTGGGGATACCATGTTCACAAGA
GAAAGACAATGAAACAGAGCCATGTGACCAA
TATGAAAGAACCATTCATGGGAAGGCCTAG
AATAGATTTTAAAACATGTTAATTCACGTTACT
ATGATTGATACCTTTAAATGTCATTTGTTGAA
ACAGACTATGATTTACAGGATCAGATGTGGA
TGTTAATTCACGTTACTTTTTTGTTAAATTTACT
AAATGTCATTTGTTGAAGGAAGATTATTCATT
CAGGATCAGATGTGGACTCTCAAATTCGACT
ATTACCTGTCAATGTTATTAATATTTTTTAGGAA
ATTCTCAGAATTTTAAACAATAACAAATCAGG

# Overview

- Genomics I (today's lecture): Focus on sequencing technology and genomes.

- Genomics II: (Friday's lecture): Focus on applications of sequencing technology.

# Overview

- Sequencing data: from wet lab to fastq.

- Applications to studying genomes and much much more.

✳Sophisticated use of data from genomics requires an integrated understanding of the biological experiment, sample preparation and down stream computational analyses of the data.
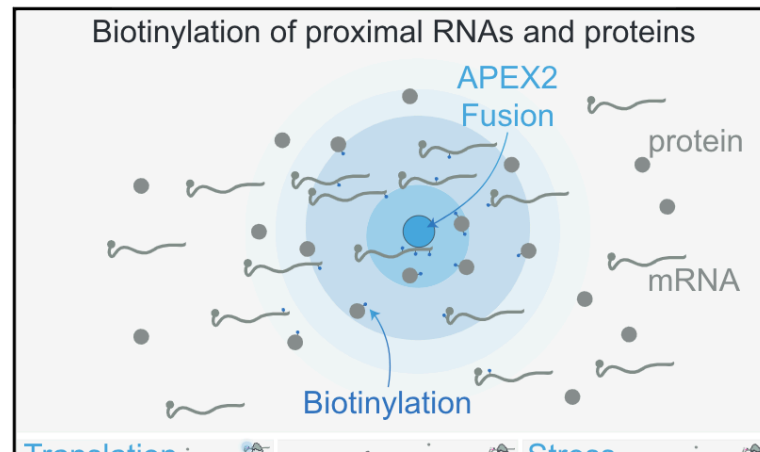
Credit: Jim Noonan for many of the slides

# Importance of genomics data: these data are central to most biomedical and biological sciences

## Molecular Cell

## Proximity RNA Labeling by APEX-Seq Reveals the Organization of Translation Initiation Complexes and Repressive RNA Granules

### Graphical Abstract



Biotinylation of proximal RNAs and proteins

APEX2 Fusion

protein

mRNA

Biotinylation

### Authors

Alejandro Padrón, Shintaro Iwasaki, Nicholas T. Ingolia

### Correspondence

ingolia@berkeley.edu

### In Brief

In this issue of *Molecular Cell*, Padrón et al. develop an RNA proximity labeling technique that maps subcellular RNA organization comprehensively. A powerful aspect of APEX-seq is the ability

**DATA AND CODE AVAILABILITY**

The raw sequencing data generated for this study are available at NCBI GEO GSE121575. Scripts to run the analyses mentioned above are available upon request.

# Data can be found in genomics databases



- Most journals require authors to submit their data to a database (e.g.,GEO) prior to publication.
- These databases entries contain raw data and processed data.
- These data can be use to examine the authors' claims, but also to test new hypotheses.

# Central questions

- Where do these data come from?

- How does the way we collect it influence what we know?

# Workflow

## 1. Isolation of sample.

*e.g.*, Isolate DNA and shear.

## 2. Library preparation

*e.g.*, Add known sequences to the ends.

# Metrics for evaluating sequencing technology

- **Throughput:**
  - Number of high quality bases per unit time
  - Number of independent samples run in parallel
  - Difficulty of sample preparation

- **Yield**
  - Number of useful reads per sample
  - Read length

- **Cost**
  - Per run cost
  - Per base cost
  - Equipment
  - Reagents
  - Labor
  - Analysis

# What is sequencing?

1. First generation sequencing

      a.  Maxam-Gilbert Sequencing
      b.  Sangar Sequencing

**2. Second generation sequencing**
      **a. Illumina Sequencing**
      b. Ion Torrent

3. Third generation sequencing

      a.  Nanopore based
      b.  Pacific Bioscience Sequencing

The technology will change, but your need to critically understand the input and output will not.

# The steps of sequencing experiments

1. Sample preparation

    a. Isolation
    b. Library construction

2. **Sequencing**

    a. Flow cell loading
    b. Cluster generation
    c. Sequencing
    d. Processing image files
    e. De-multiplexing samples

3. Data analysis

    a. Read filtering
    b. Alignment to a genome
    c. Diverse analyses



http://ycga.yale.edu/sequencing/illumina/

# What is the output from an Illumina sequencing experiment?

## One read (fastq format)

```
@HWI-ST1239:178:H0KPNADXX:2:1101:3120:1979 1:N:0:TGACCA
NCTGTAGGCTGCGTAGCCTCCCTGCAGGGTAAGTGGGAGGAGAGAGAGCAGAGGGACTTAGTGGGGCTCCCCAGGG
+
#1=DDFFFHHHHHIJIJJJJIJJJJJIJJJ?FHIDGIJ=GIHGIIIHGIJIHEHIHHGFFFFEEEDDDDDDDDDDDD
```

1. Read identifier
2. **Sequence**
3. Quality score identifier "+"
4. Quality score

# What is the output from an Illumina sequencing experiment?

## Many reads…

```
@HWI-D00306:498:HBB89ADXX:1:1101:1180:1882 1:N:0:CGATGT
NCATCACTTTCTGCACCAGCCATGACGTCAATCTTCGTCCGAACCCCAAACTCGAGATCGGAAGAGCACACGTCTG
+
#11BBDDDFFDFBFFFIIIIIIIIIIIIIFEGIIIIIFIGAGIIFIII=FEEEEEFFFDDD=@9A@BBBBB=?BB<
```

```
@HWI-D00306:498:HBB89ADXX:1:1101:1167:1902 1:N:0:CGATGT
TATTGCAATATGTTAACAATCTAACAAGGAAAAAATACCCCACACAAAACAAAACACAACCCTTAGAACTGTGCTG
+
B@@FFDFFHFHHHJJIJIGIIJJJJJJIJJHFIJJJJJIJJJEHHJJIJJJJJIIJJJJJJGHHHHFBDFFFE>CEEC
@HWI-D00306:498:HBB89ADXX:1:1101:1190:1928 1:N:0:CGATGT
ACCAAGCCACAATAAGTTAGTGTTTCCATAGTACATGCTGAGTTATTTGATCCCGTATCTATACACTGCTACTGTC
+
@<@DDDDD8CDDDGE?2<AFFBCCEEHEIEGHIIEGEIDD@CDGFFFEFIDGCFCDABFG>FBFGFGIEIFFFDDD
@HWI-D00306:498:HBB89ADXX:1:1101:1157:1931 1:N:0:CGATGT
CTGAGATTCTTTGCCATAGTCCTTAACCACTACGCAACTGCAACCAACCACCTTCCGTGGTTTGCCCTCTCGATCG
+
CCCFFFFFHHHHHHIJJIIJJJIIGHHIJGGJIGIJJJJJJJIJIIIJJJIIJJJJIIJGJJHCHFBDFFFDDECB
```

## Generally ~ 2,000,000,000 reads/sequencing lane

Note: This is for an Illumina NovaSeq with current chemistry, but this number changes

# How long are the reads?

TATTGCAATATGTTAACAATCTAACAAGGAAAAAATACCCCACACAAAACAAAACACAACCCTTAGAACTGTGCTG

←——————————————————————————————————————————————————→

**75 nt**

While there are other technologies that can give longer read lengths, Illumina reads are generally 50 nt - 250 nt

# Where do these reads come from?



Library Preparation
~2 h [15 min hands-on (Nextera)]
< 6 h [< 3 h hands-on (TruSeq)]

Cluster Generation
~5 h (<10 min hands-on)

Sequencing by Synthesis
~1.5 to 11 days

CASAVA
2 days (30 min hands-on)

Flow cell

Flow Cell A

Flow Cell B

FC A
OFF     1     2

FC B
OFF     1     2

Flow Cell Lever A

Flow Cell Lever B

# What is a flow cell?

A flow cell is a thick glass slide with 8 channels or lanes.

Each lane is randomly coated with a lawn of oligos that are complementary to library adapters

P5 oligo →

P7 oligo →

Adapters

# *Optional:* How do you make a sequencing library?

Index = unique sequence
key to identify library

P5
Rd1 SP
T
Index
Rd2 SP (P)
P7

+

(P)
A
A
(P)

Ligate

P5
Rd1 SP
Index
Rd2 SP

P7

P7
Rd2 SP
Index
Rd1 SP
P5

Amplify to create final library

P5    Rd1 SP       DNA Insert              Index
5'
                                    Rd2 SP'       P7'
                                                      5'

ends polishing

↓ 1.8X SPRI clean-up

A-tailing

↓ 1.8X SPRI clean-up

adapter ligation

↓ 1.6X SPRI clean-up

real-time PCR
amplification

↓ 1.0X SPRI clean-up

quality control
and sequencing

12 samples per lane

**Potential sources of bias:**
1. Selective PCR amplification (issue of duplicates).
2. Size selection.
3. Enzyme specificities.

Challenging but possible to analyze pg
quantities of DNA. (In humans, ~6 pg
DNA/cell).

# Where do these reads come from?



Library Preparation
~2 h [15 min hands-on (Nextera)]
< 6 h [< 3 h hands-on (TruSeq)]

Cluster Generation
~5 h (<10 min hands-on)

Sequencing by Synthesis
~1.5 to 11 days

CASAVA
2 days (30 min hands-on)



Flow cell

Flow Cell A

Flow Cell B

FC A

FC B

OFF     1     2

OFF     1     2

Flow Cell Lever A

Flow Cell Lever B

# What is the output from an Illumina sequencing experiment?

## Paired read (fastq format)

```
@HWI-ST1239:178:H0KPNADXX:2:1101:3120:1979 1:N:0:TGACCA
NCTGTAGGCTGCGTAGCCTCCCTGCAGGGTAAGTGGGAGGAGAGAGAGCAGAGGGACTTAGTGGGGCTCCCCAGGG
+
#1=DDFFFHHHHHIJIJJJIJJJJJIJJJ?FHIDGIJ=GIHGIIIHGIJIHEHIHHGFFFFEEEDDDDDDDDDDDDD

@HWI-ST1239:178:H0KPNADXX:2:1101:3120:1979 2:N:0:TGACCA
NNACCTAGCCATCTGCAGTCCTCGGTCCTGTGTTAGACCAGAACTAGGTGCCCAGGCCAGGTACCACCTAATCCTT
+
##4<@@@@@@@@@?@@@?@@??????@@??@???????????????????>?????????@>???@@@?@@??????
```

1. Read identifier
   a. Instrument
   b. Flow cell
   c. Read ID
   d. Coordinates
   e. Which read from a paired end sample
   f. Which index for multiplexed read
2. Sequence
3. Quality score identifier "+"
4. Quality score

# What limits the insert size and read length?

## One read (fastq format)

```
@HWI-D00306:498:HBB89ADXX:1:1101:1180:1882 1:N:0:CGATGT
NCATCACTTTCTGCACCAGCCATGACGTCAATCTTCGTCCGAACCCCAAACTCGAGATCGGAAGAGCACACGTCTG
+
#11BBDDDFFDFBFFFIIIIIIIIIIIIIIFEGIIIIIFIGAGIIFIII=FEEEEEFFFDDD=@9A@BBBBB=?BB<
```

- For each single end read: Incomplete incorporation of bases.

- For the size of the insert (especially for paired end analysis): Ability to get consistent clusters.

# What do I do with my sequencing reads?

# Many reference genomes are available

# There is a wide range of genome sizes.

kb = 1000 bp
Mb = $1 \times 10^6$ bp
Gb = $1 \times 10^9$ bp
Tb = $1 \times 10^{12}$ bp

# Human haploid genome ~ 3 Gb

75 nt x $3 \times 10^8$ reads/lane is about the right scale, but the amount of **coverage** necessary depends on application.



Mammals
Birds
Reptiles
Frogs
Salamanders
Lungfishes
Teleost fishes
Chondrostean fishes
Cartilaginous fishes
Jawless fishes
Non-vertebrate chordates
Crustaceans
Insects
Arachnids
Myriapods
Molluscs
Annelids
Echinoderms
Water bears (Tardigrada)
Flatworms (Platyhelminthes)
Rotifers
Red algae (Rhodophyta)
Green algae (Chlorophyta)
Brown algae (Phaeophyta)
Flowering plants (Angiosperms)
Non-flowering seed plants (Gymnosperms)
Ferns (Monilophytes)
Club mosses (Lycophytes)
Mosses and kin (Bryophytes)
Roundworms (Nematoda)
Cnidarians
Sponges (Porifera)
Fungi
Protozoa
Bacteria
Archaea

−1  0  1  2  3  4  5  6
Log$_{10}$ C-value (Mb)

1 Gb   10 Gb   100 Gb

# Sequencing of the human genome

## Victory declared **2003**







- Industrialization of Sanger sequencing, library construction, sample preparation, analysis, etc.

- $3 billion total cost

- 1 Gb/month at largest centers (2005)

Novaseq 20 billion reads 2x150 bp. $1000 -> $100/genome.

# How to assemble a genome

Generate reads

Find overlapping reads

Assemble reads into contigs

contig

Join contigs into scaffolds

*mate pair*

scaffold

Join scaffolds into "finished" sequence anchored on chromosomes

AGTTGTATTATTAGAAACTGAGGGCTAAAAACTGTGCACATACACAGACACACATATTATTTTAATATAGATTTTCAATAATTGGTCTAGGATAAGGATAATATACAG

## There are various

**Assembly quality criteria:**

Accuracy: number of errors
(Human << 1/100,000 bp)

Contiguity: number of gaps
(Human: est. 357)

**Coverage:**
Average number of reads representing a particular position in the assembly

Human, Mouse, Rat:  > 20x
Chimpanzee:      ~6x
Squirrel:    ~2x

Scaffold_0: 12,865,123 – 12,965-110

Chr5: 133,876,119 – 134,876,119

# The importance of paired end reads

Paired-End Reads

Alignment to the  Reference Sequence

Read 1

Read 2

Reference

Repeats

- Increase coverage of the insert.

- Particularly helpful when one read maps to multiple places in the genome.

```
CCAAATCAAACAGTTGTATTATTAGAAACTGAGGGCTAAAAACTGTGCACATACACAGACACACATATTATTTTAATATAGATTTTCAATAATTGGTCTAGGATAAGG
AGCAAGAAGAAACAAAGACTGTTACTATGGAAAAATGAAAATAGATTTTAAAACATGTTAATTCACGTTACTTTTTGTTAAATTTACTTTTCTTCTTTCACTTCTT
AATAAATCACATTAATTCCTTATCTCATGTGAAATTTCATATTTATGATTGATACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAATATTT
CAGTATTATGTTCTAGGCATTGGGGATACCATGTTCACAAGACAGACTATGATTTACAGGATCAGATGTGGACTCTCAAATTCGACTGAGAATAAAACAGACACT
TAATTGATGCTAGAAAGACAATGAAACAGAGCCATGTGACCAATGAGAGAGATGAGGGTGGCAGCAGCCTGTTTTAGATAAGGTACCTGATTGGTGGGATTGG
TATGCCTTAATGATATGAAAGAACCATTCATGGGAAGGCCTAGCATTAAAAACCGTCTAGGCAGAATGAGCAGCAAGTGCAAGGGTCCTGGATAGGAATGAGC
ATGGAAAAATGAAAATAGATTTTAAAACATGTTAATTCACGTTACTTTTTGTTAAATTTACTTTTCTTCTTTCACTTCTTACCTGTCAATGTTATTAATATTTTTAGGAA
GAAATTTCATATTTATGATTGATACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAATATTTTTTAGAATAATAAGTCCCAGGCACAAGA
CATGTTCACAAGACAGACTATGATTTACAGGATCAGATGTGGACTCTCAAATTCGACTGAGAATAAAACAGACACTAAACAAGTAAATAAAGTTAATTTCAAGTT
AGATTTTAAAACATGTTAATTCACGTTACTTTTTGTTAAATTTACTTTTCTTCTTTCACTTCTTACCTGTCAATGTTATTAATATTTTTAGGAACAATAAATCACATTAAT
ATTGATACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAATATTTTTTAGAATAATAAGTCCCAGGCACAAGACCAGTATTATGTTCTAG
ACTATGATTTACAGGATCAGATGTGGACTCTCAAATTCGACTGAGAATAAAACAGACACAAACAAGTAAATAAAGTTAATTTCAAGTTGTAATTGATGCTATCCA
TTGGGGATACCATTACCTGTCAATGTTATTAATATTTTTAGGAACAATAAATCACATTAATTCCAACATGCAAAGAGGAAATCTCCATATCATGCTTGTCATTCGTTT
GTGTGTAAACATTCTCAGAATTTTAAACAATAACAAATCAGGGCTGAATGTGGCCAACATGCAAAGAGGAAATCTCCCATCTGTCCAAATCAAACAGTTGTATT
CATACACAGACACACATATTATTTTAATATAGATTTTCAATAATTGGTCTAGGATAAGGATAATATACAGAGAACATGCCAAAAGTTTAAGCAAGAAGAAACAAAG
TTAAAACATGTTAATTCACGTTACTTTTTGTTAAATTTACTTTTCTTCTTTCACTTCTTACCTGTCAATGTTATTAATATTTTTAGGAACAATAAATCACATTAATTCCTT
TACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAATATTTTTTAGAATAATAAGTCCCAGGCACAAGACCAGTATTATGTTCTAGGCAT
TGATTTACAGGATCAGATGTGGACTCTCAAATTCGACTGAGAATAAAACAGACACTAAACAAGTAAATAAAGTTAATTTCAAGTTGTAATTGATGCTAGAAAGACA
AGATGAGGGTGGCAGCAGCCTGTTTTAGATAAGGTACCTGATTGGTGGGATTGGAAGACCTCTCTGAGATTAGTGTCTTCAGATATGCCTTAATGATATGAAAG
AACCGTCTAGGCAGAATGAGCAGCAAGTGCAAGGGTCCTGGATAGGAATGAGCTGGATATACTCAAGGAAGAAAGAGAAACTATGGAAAAATGAAAATAGATT
TAAATTTACTTTTCTTCTTTCACTTCTTACCTGTCAATGTTATTAATATTTTTAGGAACAATAAATCACATTAATTCCTTATCTCATGTGAAATTTCATATTTATGATTGA
TTATTCATTTTTTCATTCAATAAATATTTTTTAGAATAATAAGTCCCAGGCACAAGACCAGTATTATGTTCTAGGCATTGGGGATACCATGTTCACAAGACAGACTAT
ATTCGACTGAGAATAAAACAGACACTAAACAAGTAAATAAAGTTAATTTCAAGTTGTAATTGATGCTACTATGGAAAAATGAAAATAGATTTTAAAACATGTTAATTC
TTTCACTTCTTACCTGTCAATGTTATTAATATTTTTAGGAACAATAAATCACATTAATTCCTTATCTCATGTGAAATTTCATATTTATGATTGATACCTTTAAATGTCAT
CAATAAATATTTTTTAGAATAATAAGTCCCAGGCACAAGACCAGTATTATGTTCTAGGCATTGGGGATACCATGTTCACAAGACAGACTATGATTTACAGGATCAG
AACAGACACAAACAAGTAAATAAAGTTAATTTCAAGTTGTAATTGATGCTATCCCAGGCACAAGACCAGTATTATGTTCTAGGCATTGGGGATACCATTACCTGTC
CACATTAATTCCAACATGCAAAGAGGAAATCTCCATATCATGCTTGTCATTCGTTTATCAGAGGCCAAATGTTTTTCTTTGTAAACGTGTGTAAACATTCTCAGA
GTGGCCAACATGCAAAGAGGAAATCTCCCATCTGTCCAAATCAAACAGTTGTATTATTAGAAACTGAGGGCTAAAAACTGTGCACATACACAGACACACATTA
GATAAGGATAATATACAGAGAACATGCCAAAAGTTTAAGCAAGAAGAAACAAAGACTGTTACTATGGAAAAATGAAAATAGATTTTAAAACATGTTAATTCACGT
CTTCTTACCTGTCAATGTTATTAATATTTTTAGGAACAATAAATCACATTAATTCCTTATCTCATGTGAAATTTCATATTTATGATTGATACCTTTAAATGTCATTTGTT
AATATTTTTTAGAATAATAAGTCCCAGGCACAAGACCAGTATTATGTTCTAGGCATTGGGGATACCATGTTCACAAGACAGACTATGATTTACAGGATCAGATGTG
ACACTAAACAAGTAAATAAAGTTAATTTCAAGTTGTAATTGATGCTAGAAAGACAATGAAACAGAGCCATGTGACCAATGAGAGAGATGAGGGTGGCAGCAGC
ATTGGAAGACCTCTCTGAGATTAGTGTCTTCAGATATGCCTTAATGATATGAAAGAACCATTCATGGGAAGGCCTAGCATTAAAAACCGTCTAGGCAGAATGAG
GAGCTGGATATACTCAAGGAAGAAAGAGAAACTATGGAAAAATGAAAATAGATTTTAAAACATGTTAATTCACGTTACTTTTTGTTAAATTTACTTTTCTTCTTTCAC
GGAACAATAAATCACATTAATTCCTTATCTCATGTGAAATTTCATATTTATGATTGATACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAA
AAGACCAGTATTATGTTCTAGGCATTGGGGATACCATGTTCACAAGACAGACTATGATTTACAGGATCAGATGTGGACTCTCAAATTCGACTGAGAATAAAACAG
AGTTGTAATTGATGCTACTATGGAAAAATGAAAATAGATTTTAAAACATGTTAATTCACGTTACTTTTTGTTAAATTTACTTTTCTTCTTTCACTTCTTACCTGTCAAT
ATTAATTCCTTATCTCATGTGAAATTTCATATTTATGATTGATACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAATATTTTTTAGAATAA
TTCTAGGCATTGGGGATACCATGTTCACAAGACAGACTATGATTTACAGGATCAGATGTGGACTCTCAAATTCGACTGAGAATAAAACAGACACAAACAAGTAA
ATCCCAGGCACAAGACCAGTATTATGTTCTAGGCATTGGGGATACCATTACCTGTCAATGTTATTAATATTTTTAGGAACAATAAATCACATTAATTCCAACATGCA
TCGTTTATCAGAGGCCAAATGTTTTTCTTTGTAAACGTGTGTAAACATTCTCAGAATTTTAAACAATAACAAATCAGGGCTGAATGTGGCCAACATGCAAAGAG
TGTATTATTAGAAACTGAGGGCTAAAAACTGTGCACATACACAGACACACATATTATTTTAATATAGATTTTCAATAATTGGTCTAGGATAAGGATAATATACAGAGA
CAAAGACTGTTACTATGGAAAAATGAAAATAGATTTTAAAACATGTTAATTCACGTTACTTTTTGTTAAATTTACTTTTCTTCTTTCACTTCTTACCTGTCAATGTTA
TTCCTTATCTCATGTGAAATTTCATATTTATGATTGATACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAATATTTTTTAGAATAATAAGT
```

# What types of annotation do we have/want?

**~3 billion bp**

```
ACAATAAATCACATTAATTCCTTATCTCATGTGAAATTTCATATTTATGATTG
ATACCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAAT
AAATATTTTTTAGAATAATAAGTCCCAGGCACAAGACCAGTATTATGTTCT
AGGCATTGGGGATACCATGTTCACAAGACAGACTATGATTTACAGGATC
AGATGTGGACTCTCAAATTCGACTGAGAATAAAACAGACACTAAACAAG
TAAATAAAGTTAATTTCAAGTTGTAATTGATGCTAGAAAGACAATGAAACA
GAGCCATGTGACCAATGAGAGAGATGAGGGTGGCAGCAGCCTGTTTTA
GATAAGGTACCTGATTGGTGGGATTGGAAGACCTCTCTGAGATTAGTGT
CTTCAGATATGCCTTAATGATATGAAAGAACCATTCATGGGAAGGCCTAG
CATTAAAAACCGTCTAGGCAGAATGAGCAGCAAGTGCAAGGGTCCTGG
ATAGGAATGAGCTGGATATACTCAAGGAAGAAAGAGAAACTATGGAAAA
ATGAAAATAGATTTTAAAACATGTTAATTCACGTTACTTTTTGTTAAATTTA
CTTTTCTTCTTTCACTTCTTACCTGTCAATGTTATTAATATTTTTAGGAACA
ATAAATCACATTAATTCCTTATCTCATGTGAAATTTCATATTTATGATTGATA
CCTTTAAATGTCATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAA
TATTTTTTAGAATAATAAGTCCCAGGCACAAGACCAGTATTATGTTCTAGG
CATTGGGGATACCATGTTCACAAGACAGACTATGATTTACAGGATCAGAT
GTGGACTCTCAAATTCGACTGAGAATAAAACAGACACTAAACAAGTAAAT
AAAGTTAATTTCAAGTTGTAATTGATGCTACTATGGAAAAATGAAAATAGA
TTTTAAAACATGTTAATTCACGTTACTTTTTGTTAAATTTACTTTTCTTCTTT
CACTTCTTACCTGTCAATGTTATTAATATTTTTAGGAACAATAAATCACATT
AATTCCTTATCTCATGTGAAATTTCATATTTATGATTGATACCTTTAAATGT
CATTTGTTGAAGGAAGATTATTCATTTTTTCATTCAATAAATATTTTTTAGA
ATAATAAGTCCCAGGCACAAGACCAGTATTATGTTCTAGGCATTGGGGAT
ACCATGTTCACAAGACAGACTATGATTTACAGGATCAGATGTGGACTCTC
AAATTCGACTGAGAATAAAACAGACACAAACAAGTAAATAAAGTTAATTT
CAAGTTGTAATTGATGCTATCCCAGGCACAAGACCA....
```

## Genes:
- Coding, noncoding, miRNA, etc.
- Isoforms
- Expression

## Genetic variation:
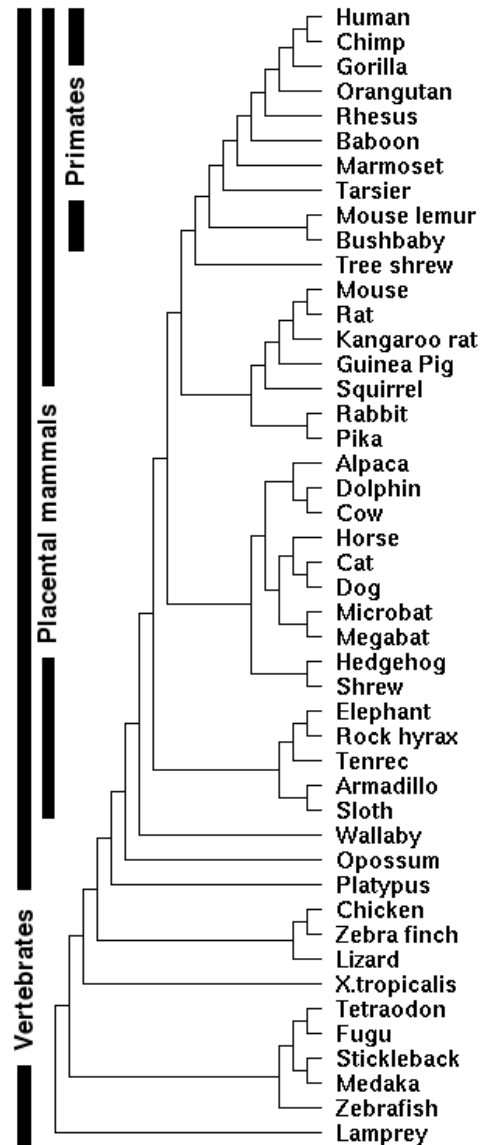- SNPs and CNVs

## Sequence conservation

## Regulatory sequences:
- Promoters
- Enhancers
- Insulators

## Epigenetics:
- DNA methylation
- Chromatin

# Degrees of genomic annotation vary widely



## ENCODE and modENCODE

**Human, Mouse (Fly, Worm, Yeast):**
- Chromosome assemblies
- Dense gene and regulatory maps, variation, etc.

**Other models (Dog, Chicken, Zebrafish):**
- Chromosome assemblies
- Partial gene maps; variation; little regulatory data

**Low coverage vertebrate genomes:**
- Scaffold assemblies
- Few annotated genes
- Used for comparative purposes

# Where do you look for existing annotations?

**UCSC Genome Browser** (genome.ucsc.edu):
   Visualization, data recovery, simple analysis
   (also http://genome-preview.ucsc.edu/)

**ENSEMBL** (ensembl.org):
   Visualization, data recovery, simple analysis

**Integrative Genomics Viewer**
(broadinstitute.orgsoftware/igv/):
   Local genome viewer (visualize local and remote data)

**Galaxy** (main.g2.bx.psu.edu):
   Complex data analysis and workflows

# Example of a genome browser track (UCSC)

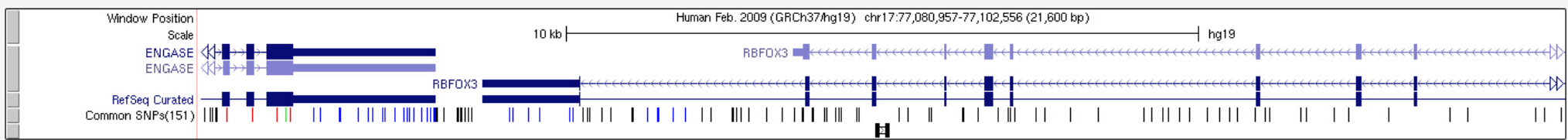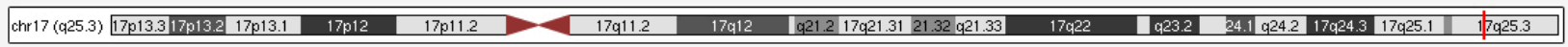Chr5: 133,876,119 – 134,876,119

# Our specific example:

```
@HWI-ST1239:178:H0KPNADXX:2:1101:3120:1979 1:N:0:TGACCA
NCTGTAGGCTGCGTAGCCTCCCTGCAGGGTAAGTGGGAGGAGAGAGAGCAGAGGGACTTAGTGGGGCTCCCCAGGG
+
#1=DDFFFHHHHHIJIJJJIJJJJJIJJJ?FHIDGIJ=GIHGIIIHGIJIHEHIHHGFFFFEEEDDDDDDDDDDDD

@HWI-ST1239:178:H0KPNADXX:2:1101:3120:1979 2:N:0:TGACCA
NNACCTAGCCATCTGCAGTCCTCGGTCCTGTGTTAGACCAGAACTAGGTGCCCAGGCCAGGTACCACCTAATCCTT
+
##4<@@@@@@@@@?@@@?@@??????@@??@?????????????????>??????????@>???@@@?@@??????
```

# Workflow

## 1. Isolation of sample.

*e.g.*, Isolate DNA and shear.

## 2. Library preparation

*e.g.*, Add known sequences to the ends.

# Using sequencing to annotate the genome

1. Where are the cis-acting regulatory elements in DNA?
   - A. DNase I hyper-sensitivity mapping (DNase-Seq).
   - B. FAIRE to map regulatory elements.

2. Where do transcription factors bind?
   - C. ChIP-seq of transcription factors (or in high res, ChIP-exo)
   - D. Nucleosome mapping (MNase-Seq).

3. Where are different histone modifications found?
   - E. ChIP-Seq of histone modifications.
   - F. ChIP-Seq of chromatin writers, readers and erasers.

4. Where is RNA polymerase transcribing?
   - G. ChIP-Seq of polymerase.
   - H. GRO-Seq, NET-Seq and TT-Seq to measure RNA in the polymerase active site..

5. How is the genome organized in 3D?
   - I. 4C/5C/Hi-C to measure chromatin conformation.

Applications of sequencing technology next week.

# Conclusions

- High-throughput sequencing has become democratized - moved out of industrial-scale genome centers

- Sequence is no longer limiting - next generation of sequencers will make sequencing very inexpensive

- Earlier methods for counting / resequencing applications are largely obsolete

- Scale of data production outstripping our ability to store and analyze it