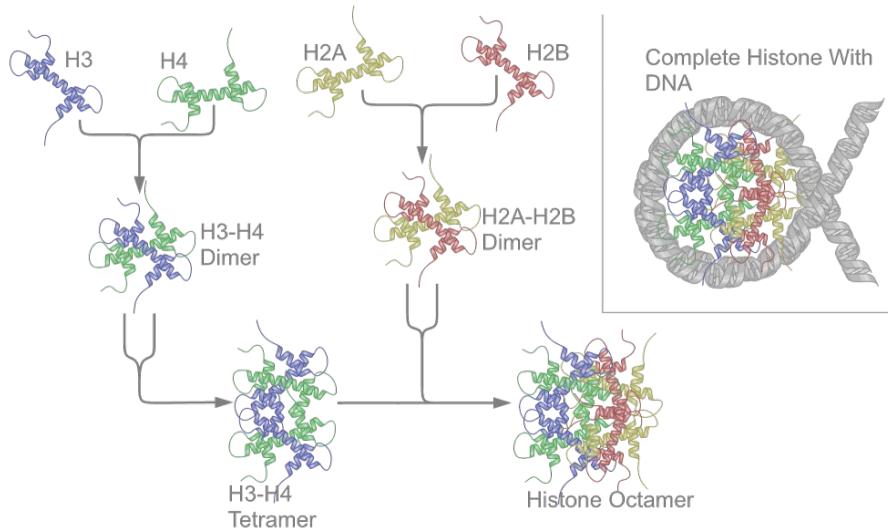


# Eukaryotic Genome Structure and Features

Genome Biology (BIOL7263)  
24Oct24

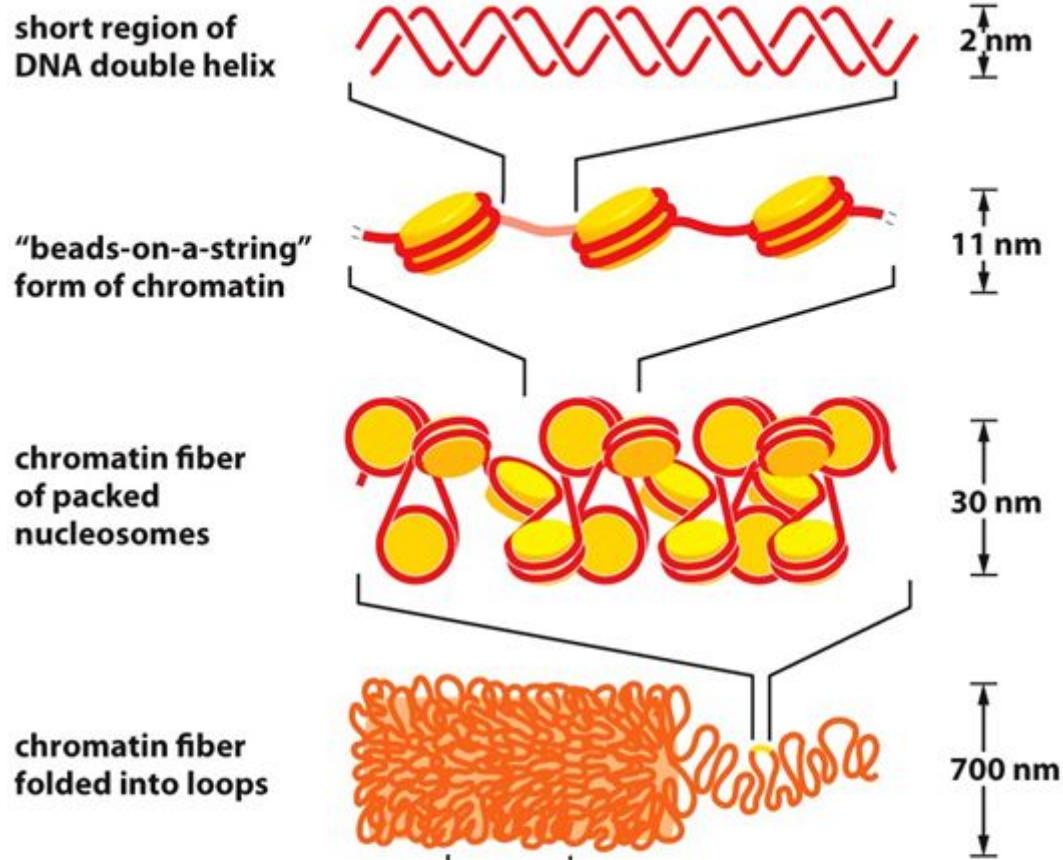
# The genomes of Eukaryotes is organized into linear chromosomes

- Chromosomes are made up of chormatin - DNA and protein complex
- DNA is wrapped around histones composed of 8 protein molecules



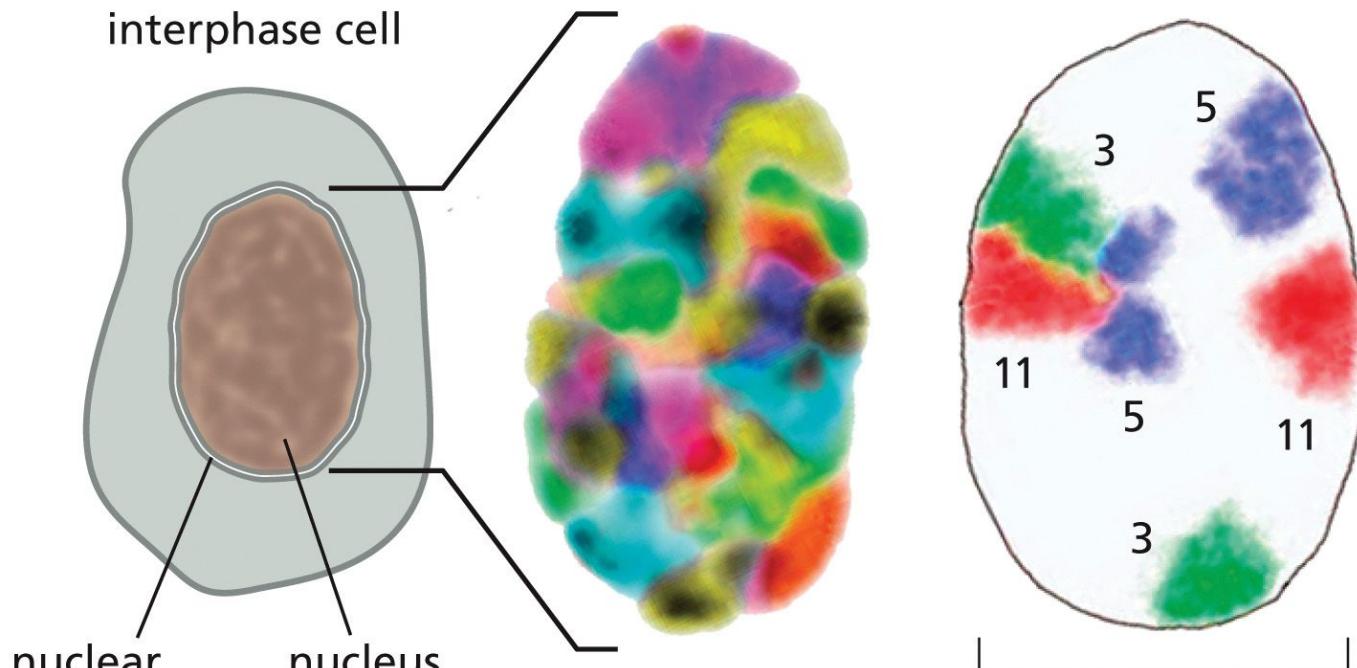
- Histone proteins are positively charged and highly conserved among organisms

These  
“beads-on-a-string”  
are further  
condensed into  
chromatin fibers

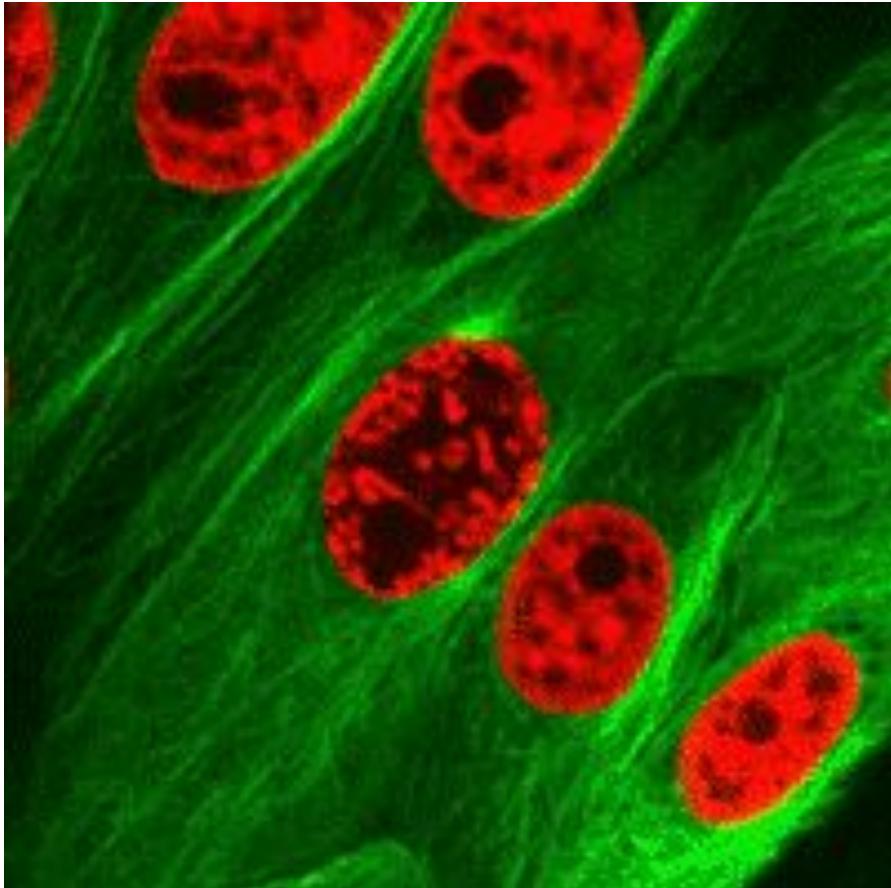


# Chromatin is dynamic

- The packing and condensation of chromatin changes with **cell cycle stage** and **transcriptional state**.



# Metaphase



- Chromosomes duplicated and highly condensed
- Essential step in the distribution of the genome to daughter cells

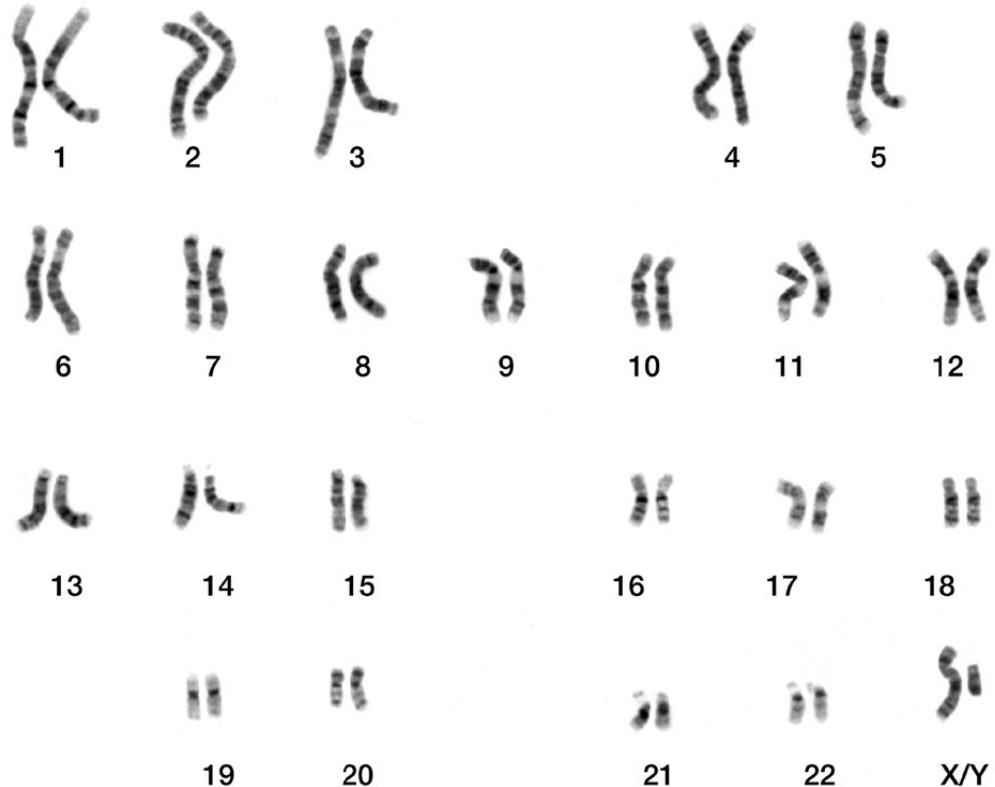
# Metaphase chromosome banding and karyogram

**TABLE 7.1 STAINING TECHNIQUES USED TO PRODUCE CHROMOSOME BANDING PATTERNS**

Technique	Procedure	Banding pattern
G-banding	Mild proteolysis followed by staining with Giemsa	Dark bands are AT-rich Pale bands are GC-rich
R-banding	Heat denaturation followed by staining with Giemsa	Dark bands are GC-rich Pale bands are AT-rich
Q-banding	Stain with quinacrine	Dark bands are AT-rich Pale bands are GC-rich
C-banding	Denature with barium hydroxide and then stain with Giemsa	Dark bands contain constitutive heterochromatin (see Section 10.1)

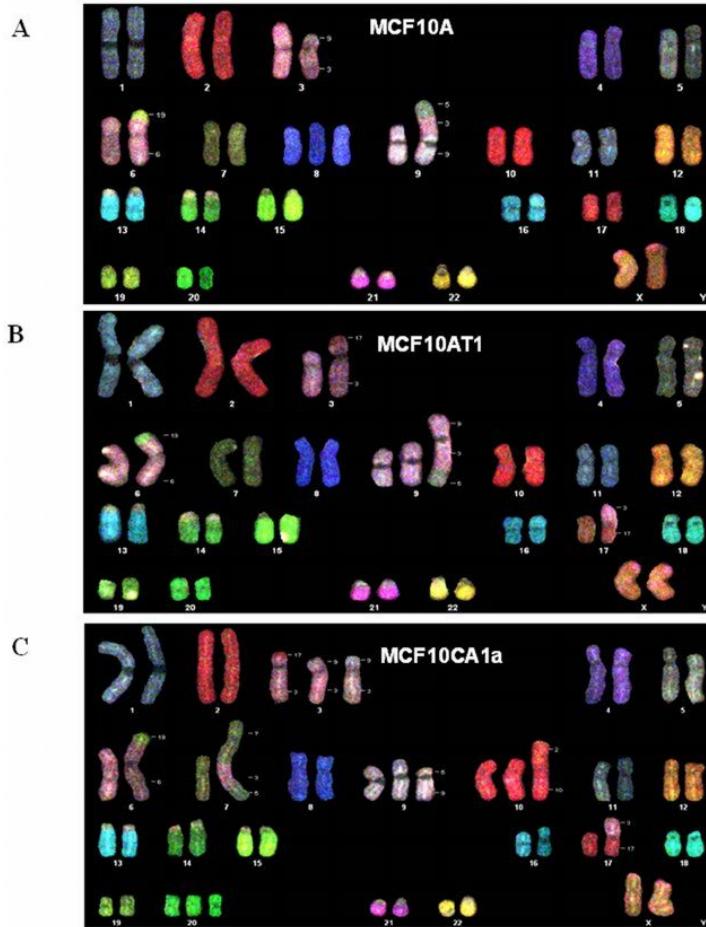
# Human Karyotype

- G-band staining



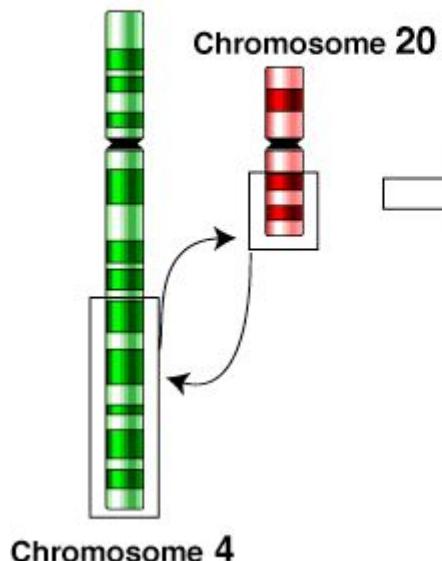
# Karyotyping is an important diagnostic tool

- Cancer progression often involves non-disjunction and abnormal chromosome numbers (Aneuploidy)

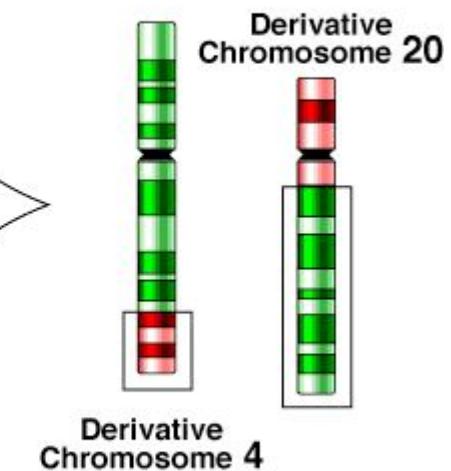


Karyotyping also can  
be used to detect  
translocation events

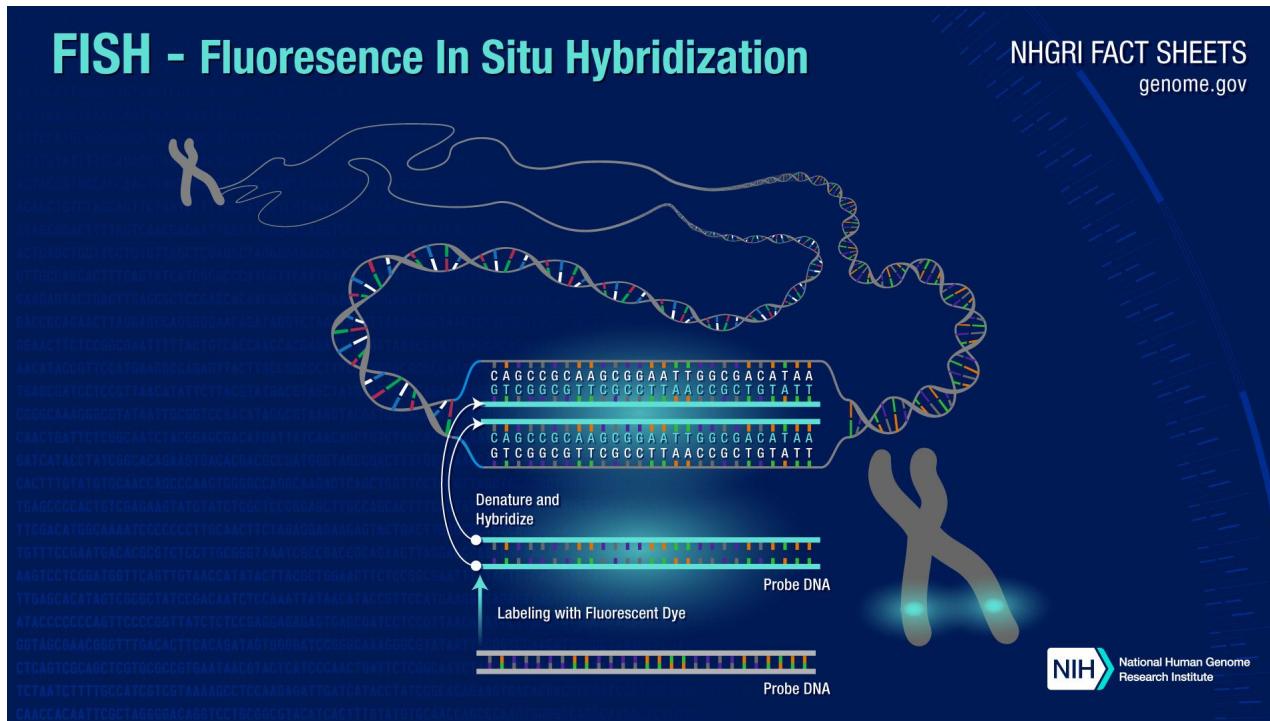
Before translocation



After translocation

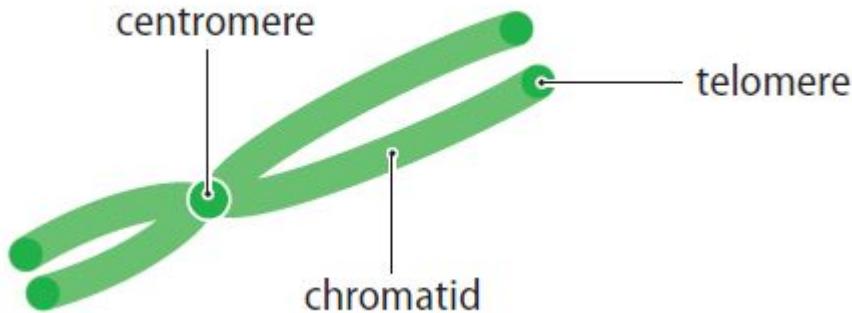


Fluorescence *in situ* hybridization (FISH) is an important tool for mapping specific sequences to chromosome structure



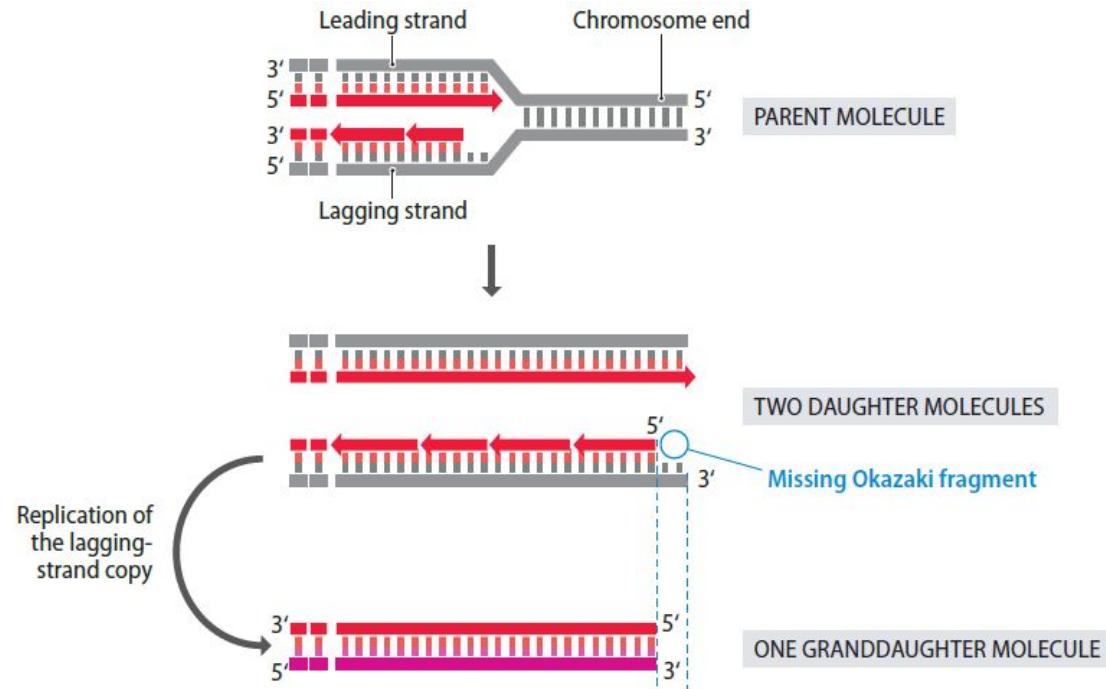
# Metaphase chromosome features

- **Chromatid** - one copy of the duplicated chromosome
- **Telomere** - Highly repetitive sequences at the end of the chromosome that protect against degradation.
- **Centromere** - specialized DNA sequences and associated proteins essential for disjunction of the duplicated chromosomes.



# Telomeres

- The replication of linear DNA molecules results in progressive shortening of the chromosomes

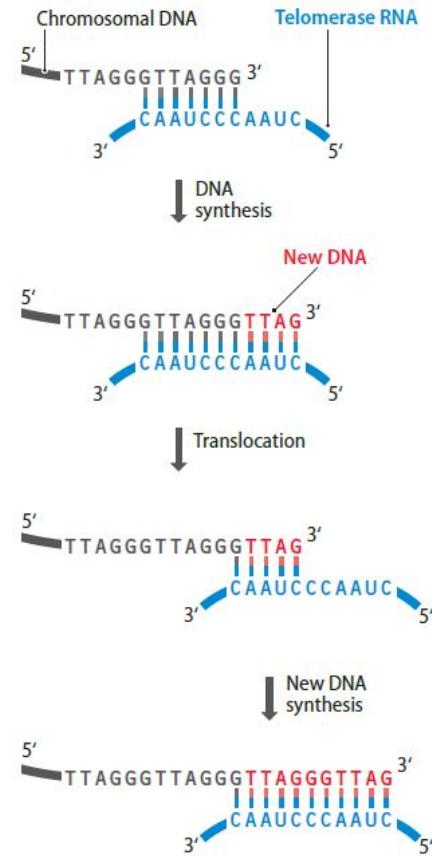


# Telomeres

- The enzyme telomerase extends the chromosome end sequence by reverse transcribing an RNA template with repeating sequence of bases

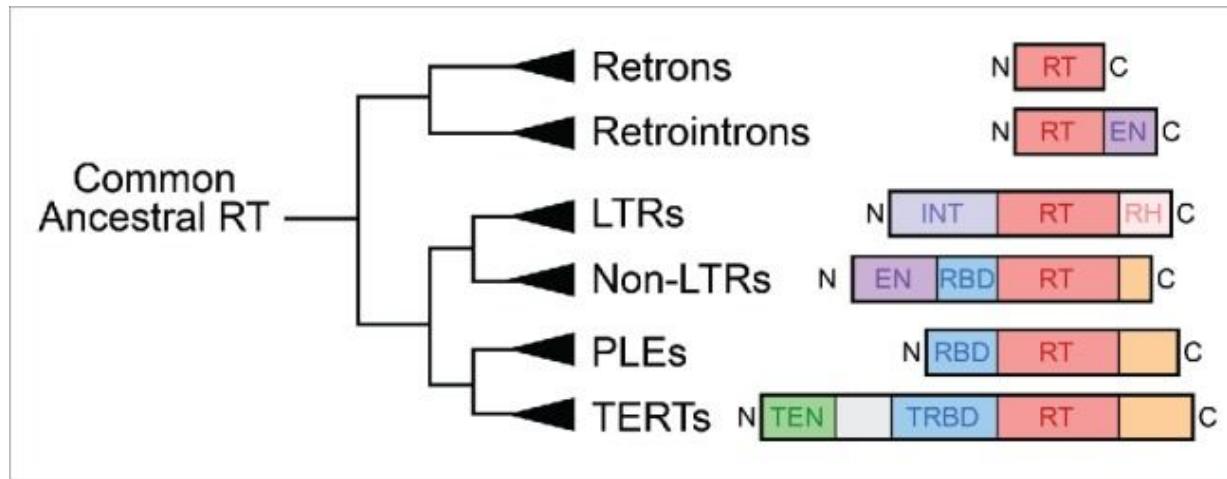
**TABLE 15.3 SEQUENCES OF TELOMERE REPEATS AND TELOMERASE RNAs IN VARIOUS ORGANISMS**

Species	Telomere repeat sequence	Telomerase RNA template sequence
Human	5'-TTAGGG-3'	5'-CUAACCCUAAC-3'
Oxytricha	5'-TTTGGGG-3'	5'-CAAAACCCAAAACC-3'
Tetrahymena	5'-TTGGGG-3'	5'-CAACCCCAA-3'



# Telomeres

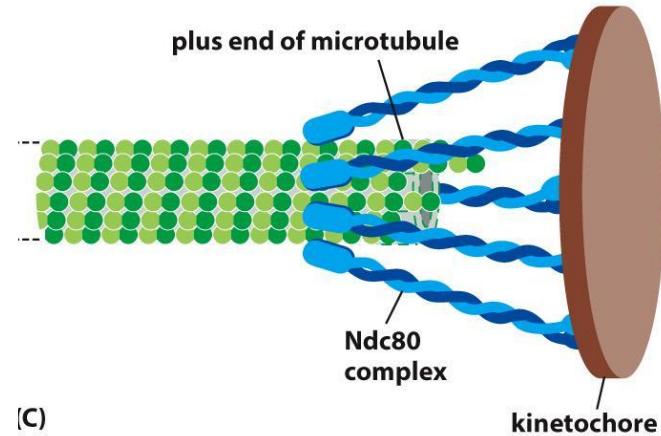
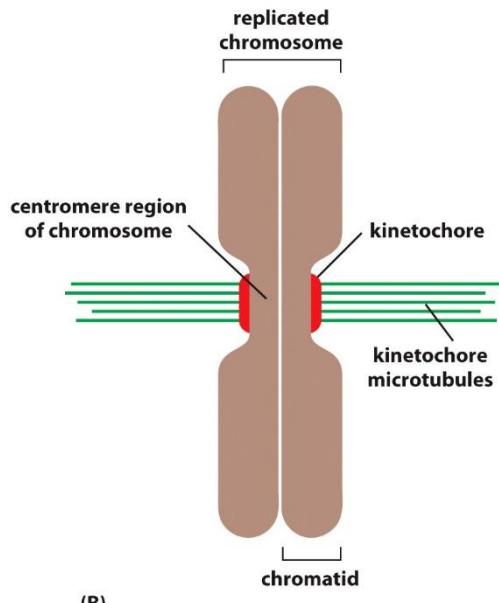
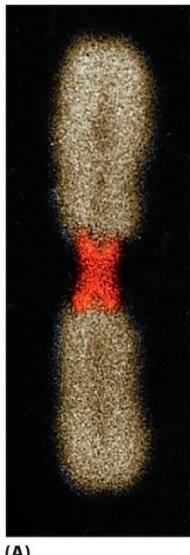
- Telomerase (TERT) may derive from reverse transcriptases carried by retrotransposons!



The central catalytic RT domain (red) is flanked by variable accessory domains, including endonuclease (EN, violet), integrase (INT, indigo), RNase-H (RH, pink), RNA binding domain (RBD, blue), and a thumb domain (orange). TERT contains a large N-terminal extension comprising of the DNA binding TEN (green) domain and TR binding domain (TRBD, blue).

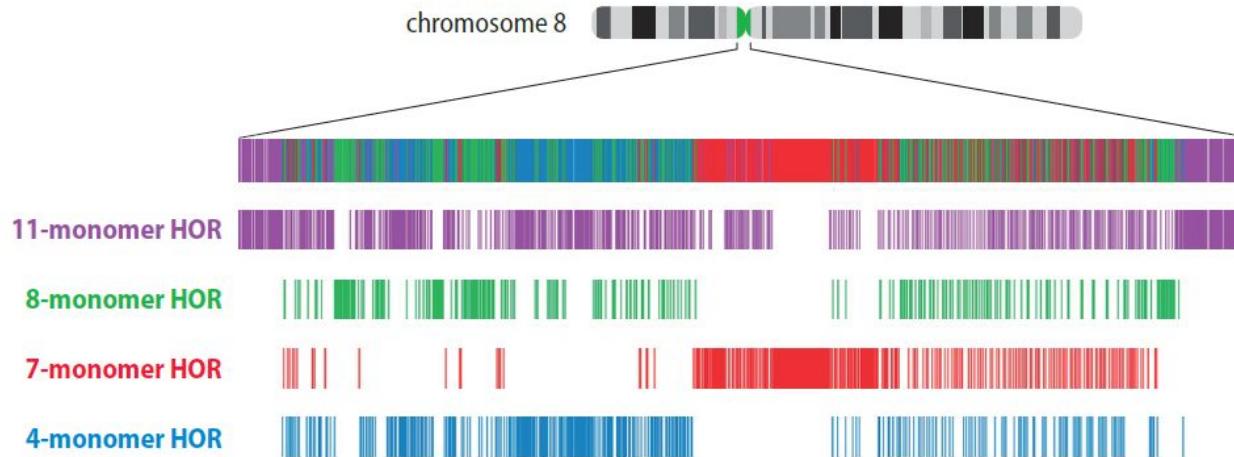
# The centromere

- Cohesin proteins bind and hold sister chromatids together
- Assembly point for the kinetochore [define]



# The centromere

- Centromere sequence is millions of bases long
- Repetitive alphoid DNA
  - 171 bp units
  - Arranged head-to-tail arrays in highly-order repeats (HORs)
  - Varying numbers of repeating units with location and among individuals
  - Impossible to resolve with short-read sequencing



# Telomere-to-telomere sequence of a human genome

## RESEARCH ARTICLE

### HUMAN GENOMICS

#### The complete sequence of a human genome

Sergey Nurk<sup>1†</sup>, Sergey Koren<sup>1†</sup>, Arang Rhie<sup>1</sup>, Mikko Rautiainen<sup>1</sup>, Andrey V. Bzikadze<sup>2</sup>, Alla Mikheenko<sup>3</sup>, Mitchell R. Vollger<sup>4</sup>, Nicolas Altemose<sup>5</sup>, Lev Uralsky<sup>6,7</sup>, Ariel Gershman<sup>8</sup>, Sergey Aganezov<sup>9‡</sup>, Savannah J. Hoyt<sup>10</sup>, Mark Diekhans<sup>11</sup>, Glennis A. Logsdon<sup>4</sup>, Michael Alonge<sup>9</sup>, Stylianos E. Antonarakis<sup>12</sup>, Matthew Borchers<sup>13</sup>, Gerard G. Bouffard<sup>14</sup>, Shelise Y. Brooks<sup>14</sup>, Gina V. Caldas<sup>15</sup>, Nae-Chyun Chen<sup>9</sup>, Haoyu Cheng<sup>16,17</sup>, Chen-Shan Chin<sup>18</sup>, William Chow<sup>19</sup>, Leonardo G. de Lima<sup>13</sup>, Philip C. Dishuck<sup>4</sup>, Richard Durbin<sup>19,20</sup>, Tatiana Dvorkina<sup>3</sup>, Ian T. Fiddes<sup>21</sup>, Giulio Formenti<sup>22,23</sup>, Robert S. Fulton<sup>24</sup>, Arkarachai Fungtammasan<sup>18</sup>, Erik Garrison<sup>11,25</sup>, Patrick G. S. Grady<sup>10</sup>, Tina A. Graves-Lindsay<sup>26</sup>, Ira M. Hall<sup>27</sup>, Nancy F. Hansen<sup>28</sup>, Gabrielle A. Hartley<sup>10</sup>, Marina Haukness<sup>11</sup>, Kerstin Howe<sup>10</sup>, Michael W. Hunkapiller<sup>29</sup>, Chirag Jain<sup>1,30</sup>, Miten Jain<sup>11</sup>, Erich D. Jarvis<sup>22,23</sup>, Peter Kerpeljiev<sup>31</sup>, Melanie Kirsche<sup>9</sup>, Mikhail Kolmogorov<sup>32</sup>, Jonas Korlach<sup>29</sup>, Milinn Kremitzki<sup>26</sup>, Heng Li<sup>16,17</sup>, Valerie V. Maduro<sup>33</sup>, Tobias Marschall<sup>34</sup>, Ann M. McCartney<sup>1</sup>, Jennifer McDaniel<sup>35</sup>, Danny E. Miller<sup>4,36</sup>, James C. Mullikin<sup>14,28</sup>, Eugene W. Myers<sup>37</sup>, Nathan D. Olson<sup>35</sup>, Benedict Paten<sup>11</sup>, Paul Peluso<sup>29</sup>, Pavel A. Pevzner<sup>32</sup>, David Purkay<sup>4</sup>, Tamara Potapova<sup>13</sup>, Evgeny I. Rogayev<sup>6,7,38,39</sup>, Jeffrey A. Rosenfeld<sup>40</sup>, Steven L. Salzberg<sup>9,41</sup>, Valerie A. Schneider<sup>42</sup>, Fritz J. Sedlazeck<sup>43</sup>, Kishwan Shafin<sup>11</sup>, Colin J. Shew<sup>44</sup>, Alaina Shumate<sup>41</sup>, Ying Sims<sup>19</sup>, Arian F. A. Smit<sup>45</sup>, Daniela C. Soto<sup>44</sup>, Ivan Sovic<sup>29,46</sup>, Jessica M. Storer<sup>45</sup>, Aaron Streets<sup>5,47</sup>, Beth A. Sullivan<sup>48</sup>, Françoise Thibaud-Nissen<sup>42</sup>, James Torrance<sup>19</sup>, Justin Wagner<sup>35</sup>, Brian P. Walenz<sup>1</sup>, Aaron Wenger<sup>20</sup>, Jonathan M. D. Wood<sup>19</sup>, Chunlin Xiao<sup>42</sup>, Stephanie M. Yan<sup>49</sup>, Alice C. Young<sup>14</sup>, Samantha Zarate<sup>9</sup>, Urvashi Surti<sup>50</sup>, Rajiv C. McCoy<sup>49</sup>, Megan Y. Dennis<sup>44</sup>, Ivan A. Alexandrov<sup>3,7,51</sup>, Jennifer L. Gerton<sup>13,52</sup>, Rachel J. O'Neill<sup>10</sup>, Winston Timp<sup>8,41</sup>, Justin M. Zook<sup>35</sup>, Michael C. Schatz<sup>9,49</sup>, Evan E. Eichler<sup>4,53\*</sup>, Karen H. Miga<sup>11,54\*</sup>, Adam M. Phillippy<sup>1\*</sup>

Since its initial release in 2000, the human reference genome has covered only the euchromatic fraction of the genome, leaving important heterochromatic regions unfinished. Addressing the remaining 8% of the genome, the Telomere-to-Telomere (T2T) Consortium presents a complete 3.055 billion-base pair sequence of a human genome, T2T-CHM13, that includes gapless assemblies for all chromosomes except Y, corrects errors in the prior references, and introduces nearly 200 million base pairs of sequence containing 1956 gene predictions, 99 of which are predicted to be protein coding. The completed regions include all centromeric satellite arrays, recent segmental duplications, and the short arms of all five acrocentric chromosomes, unlocking these complex regions of the genome to variational and functional studies.

T2T-CHM13 fills 8% of the reference which was incomplete

- ~200 Mbps new sequence
- 99 new protein coding genes
- Resolution of telomere and centromere sequences of each chromosome

Nurk, S., S. Koren, A. Rhie, M. Rautiainen, A. V. Bzikadze, A. Mikheenko, M. R. Vollger, N. Altemose, L. Uralsky, A. Gershman, S. Aganezov, et al. (2022). The complete sequence of a human genome. *Science* (New York, N.Y.) 376:44–53.

# Telomere-to-telomere sequence of a human genome

One weird trick to sequencing the entire genome:

## Complete hydatidiform mole (CHM)

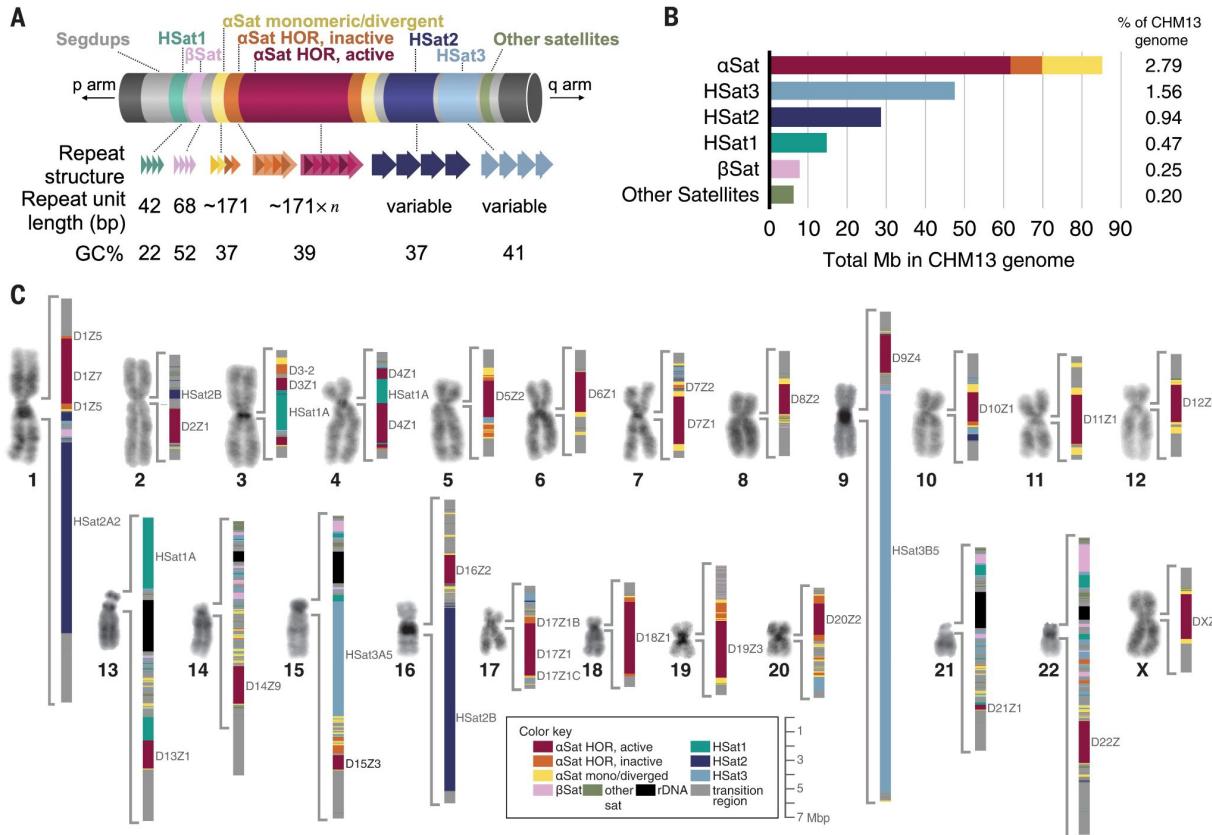
- non-viable fertilized egg, missing the maternal nucleus
- duplication of the paternal complement post-fertilization produces a perfectly homozygous genome.

Then:

- 30x PacBio HiFi
- 120x Nanopore ultralong reads
- 100x illumina
- 70x illumina HiC

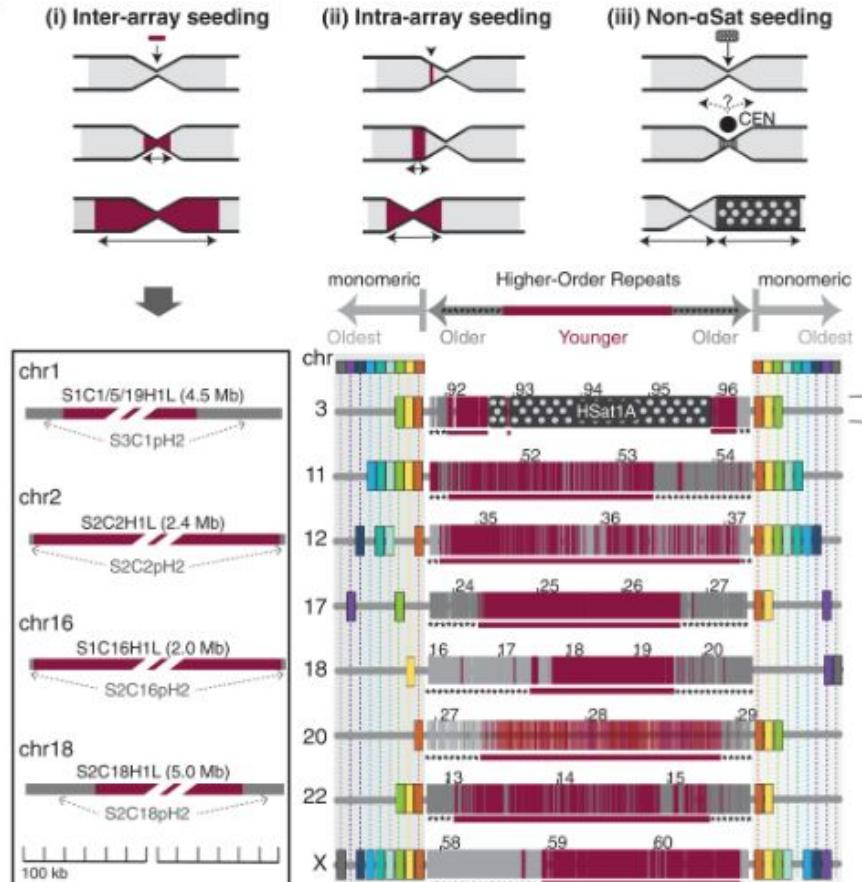
# Telomere-to-telomere sequence of a human genome

- Resolves sequence of the centromere
- Reveals centromeric sequences are expanding through tandem duplication



# Telomere-to-telomere sequence of a human genome

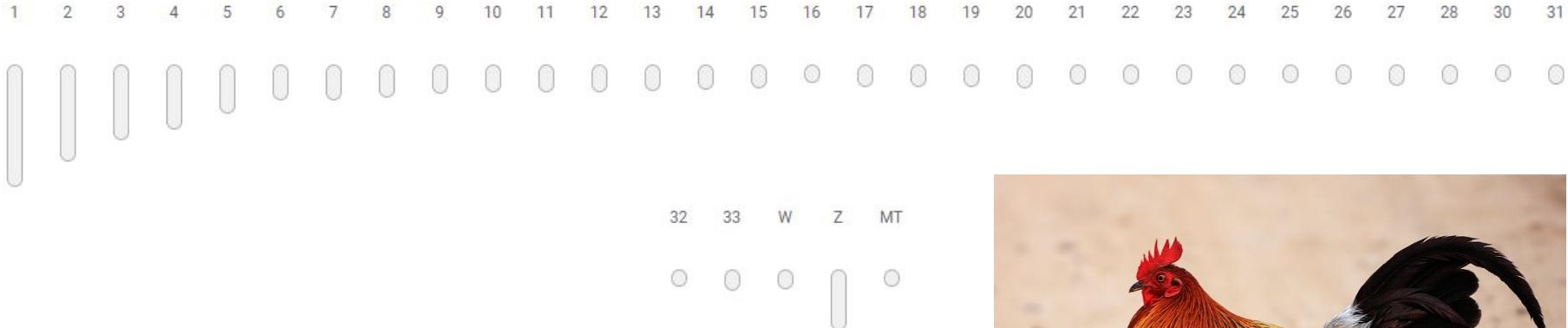
- Resolves sequence of the centromere
- Reveals centromeric sequences are expanding through tandem duplication
- Youngest sequences are the ones interacting with kinetochore proteins



Altemose, N., G. A. Logsdon, A. V. Bzikadze, P. Sidhwani, S. A. Langley, G. V. Caldas, S. J. Hoyt, L. Uralsky, F. D. Ryabov, C. J. Shew, M. E. G. Sauria, et al. (2022). Complete genomic and epigenetic maps of human centromeres. *Science* (New York, N.Y.) 376:eabl4178.

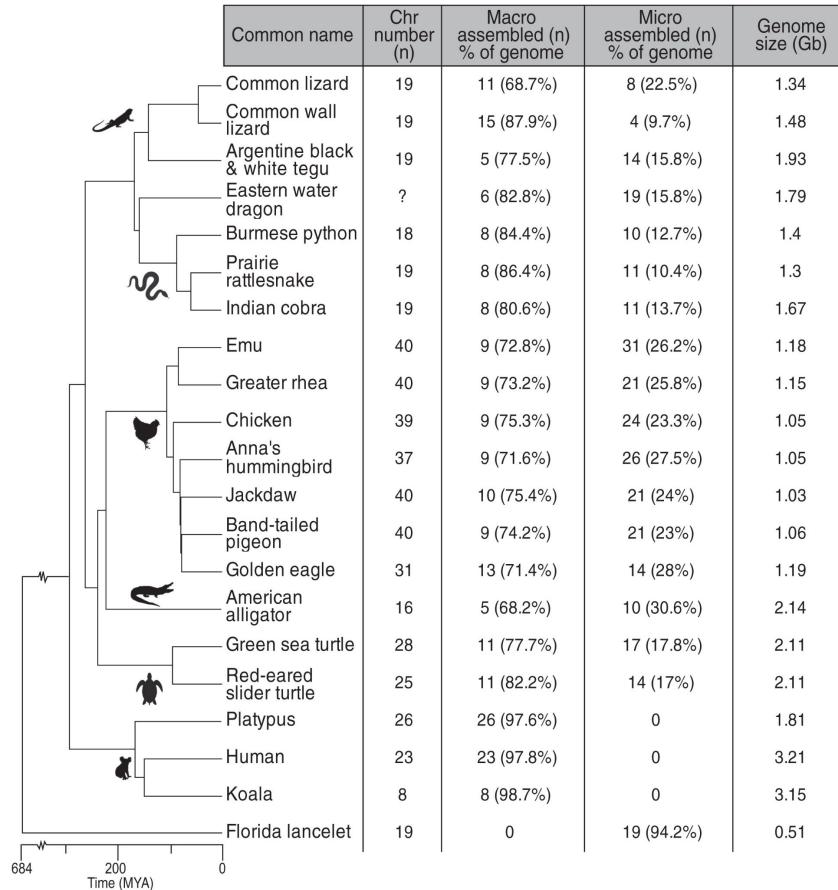
# Chromosomal Diversity

- Microchromosomes - < 20 Mb in length - birds, fish, reptiles and amphibians
- Gene-dense, low repetitive sequence content, and have high rates of recombination.



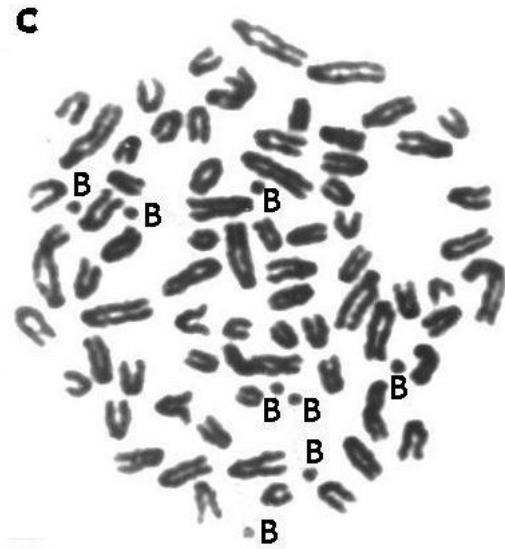
# Chromosomal Diversity

- Microchromosomes may be the ancestral state and macrochromosomes a derived state in mammals
- How and why these changes have occurred remains unresolved



# Chromosomal Diversity

- B chromosomes
  - Non-essential, atypical chromosomes found in some individuals
  - Usually fragments of other chromosomes and may contain genes including rRNAs

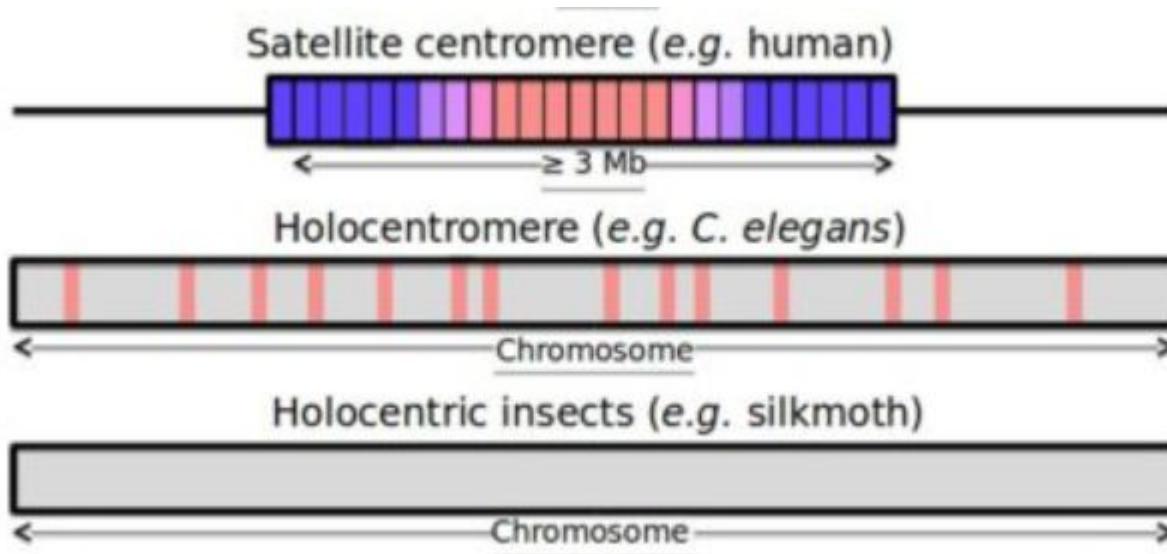


Roe deer karyotype

<https://molecularcytogenetics.biomedcentral.com/articles/10.1186/1755-8166-4-22>

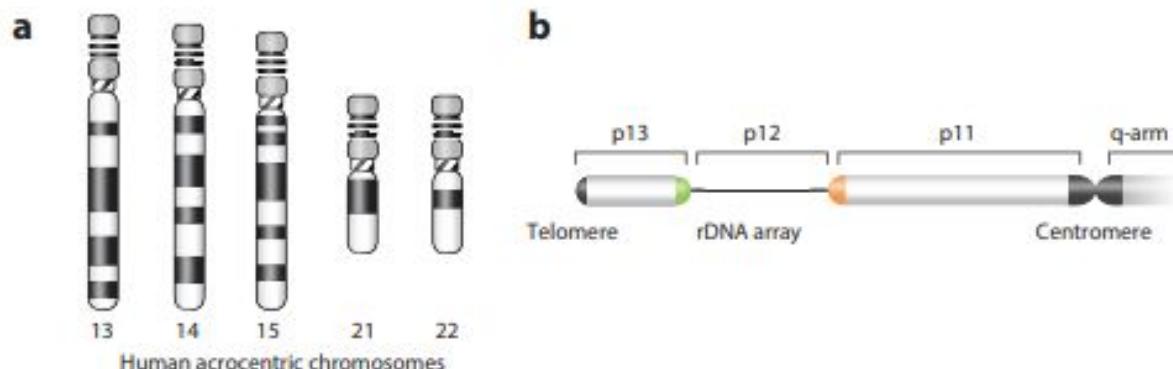
# Chromosomal Diversity

- Holocentric chromosomes - *C. elegans*, some insects, some plants - multiple centromere-like structures distributed along the chromosomes or entirely absent.



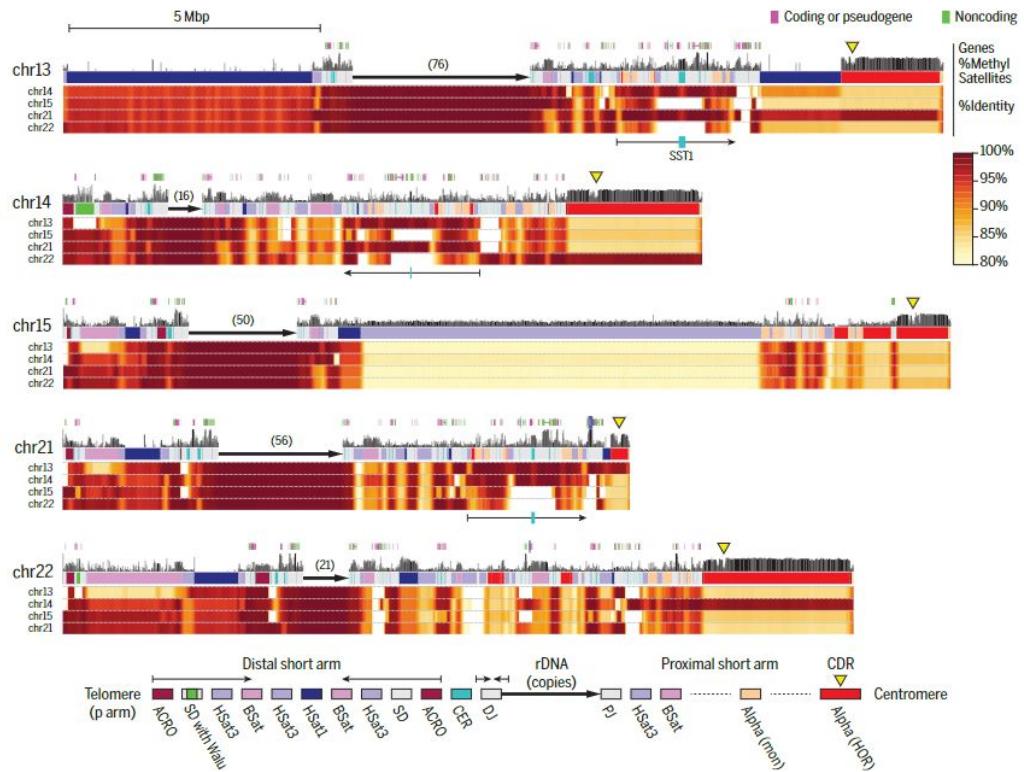
# Chromosomal Diversity

- Acrocentric chromosomes - chromosomes with centromere located close to the end of the chromosome
- In humans the short arm sequences are highly repetitive and high rates of heterologous recombination among chromosomes



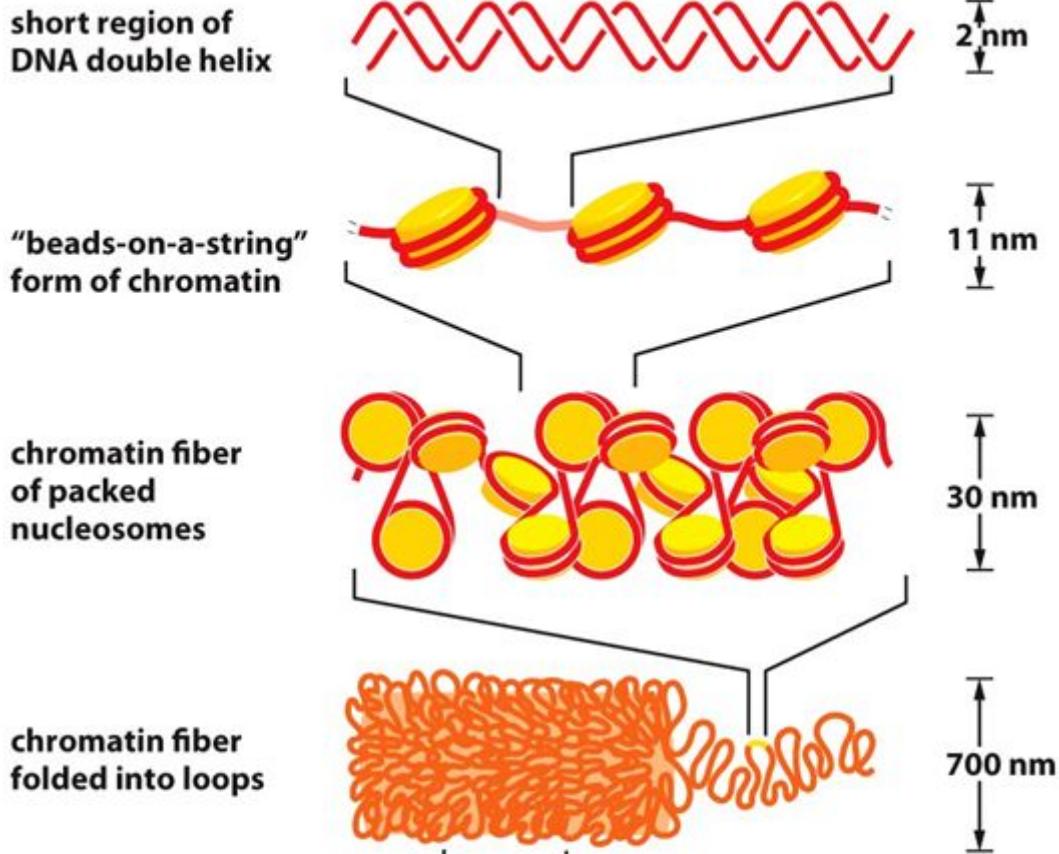
# Chromosomal Diversity

- In humans the short arm sequences are highly repetitive and high rates of heterologous recombination among chromosomes
- rDNA copy numbers highly variable and among populations and individuals



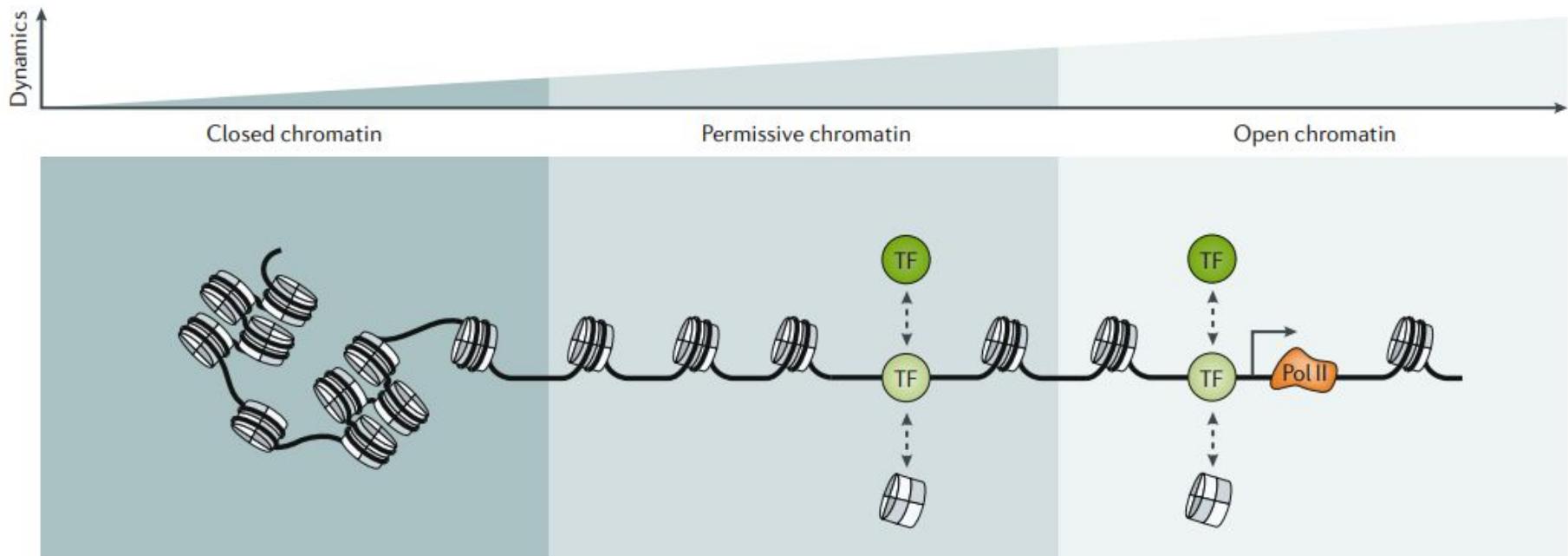
# Chromatin state

- Chromatin must be “opened” for gene expression to proceed



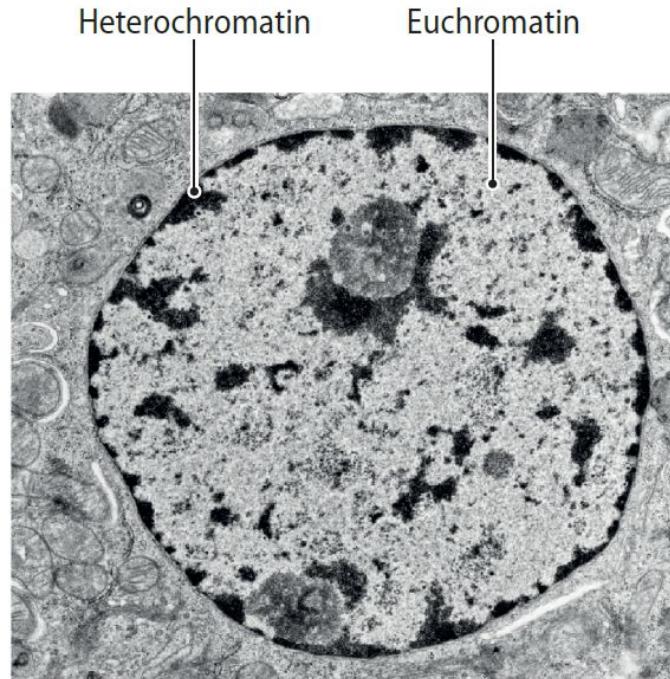
# Chromatin state

- Chromatin must be “opened” for gene expression to proceed



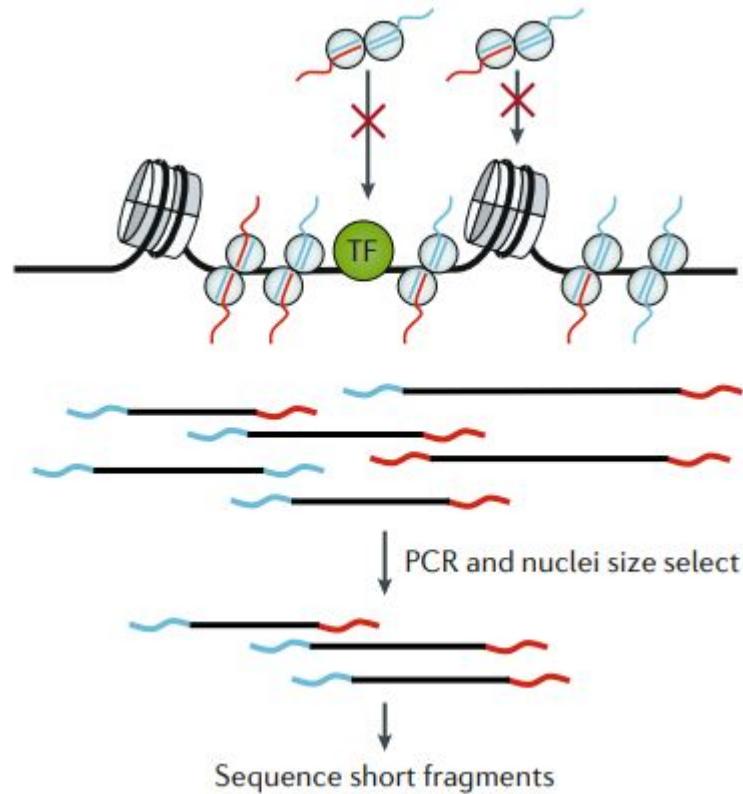
# Chromatin state

- Tightly condensed **heterochromatin** in transcriptionally inactive
- Less dense **euchromatin** is the site of active transcription

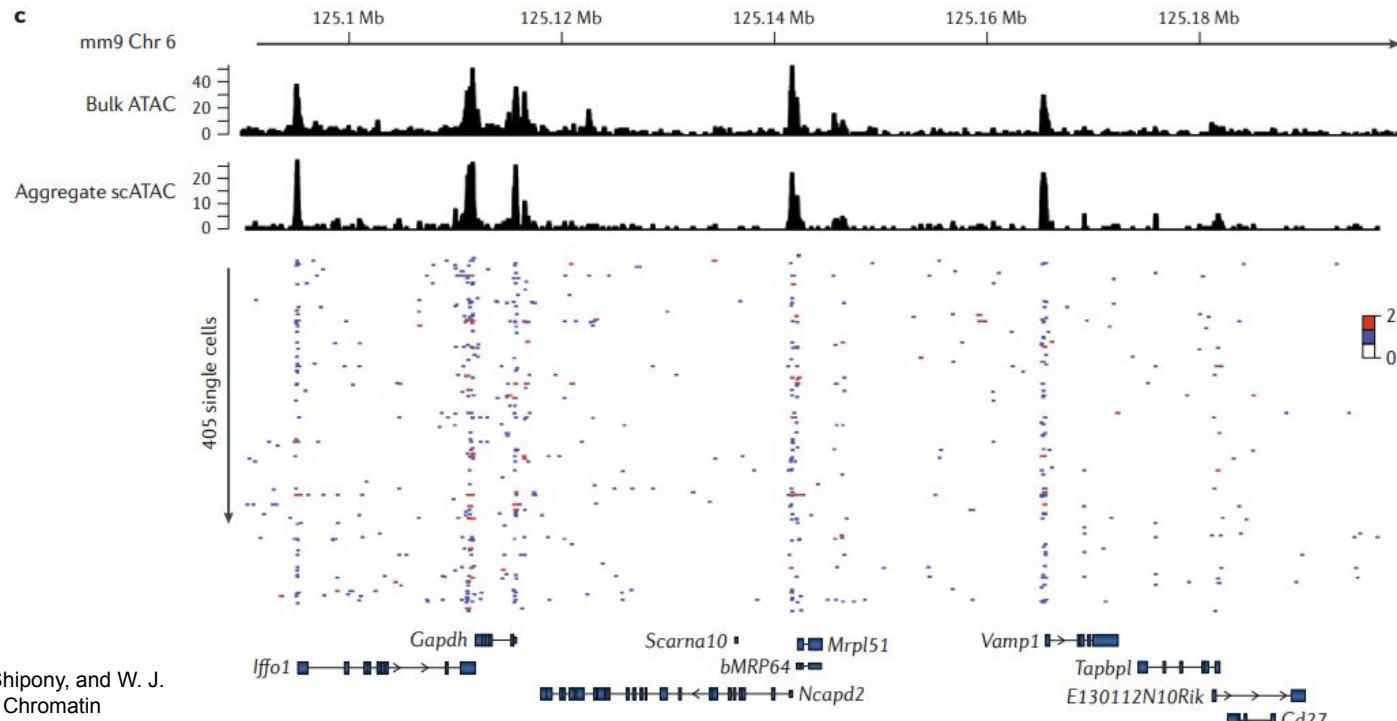


# Mapping open chromatin can yield insights into transcriptional regulation

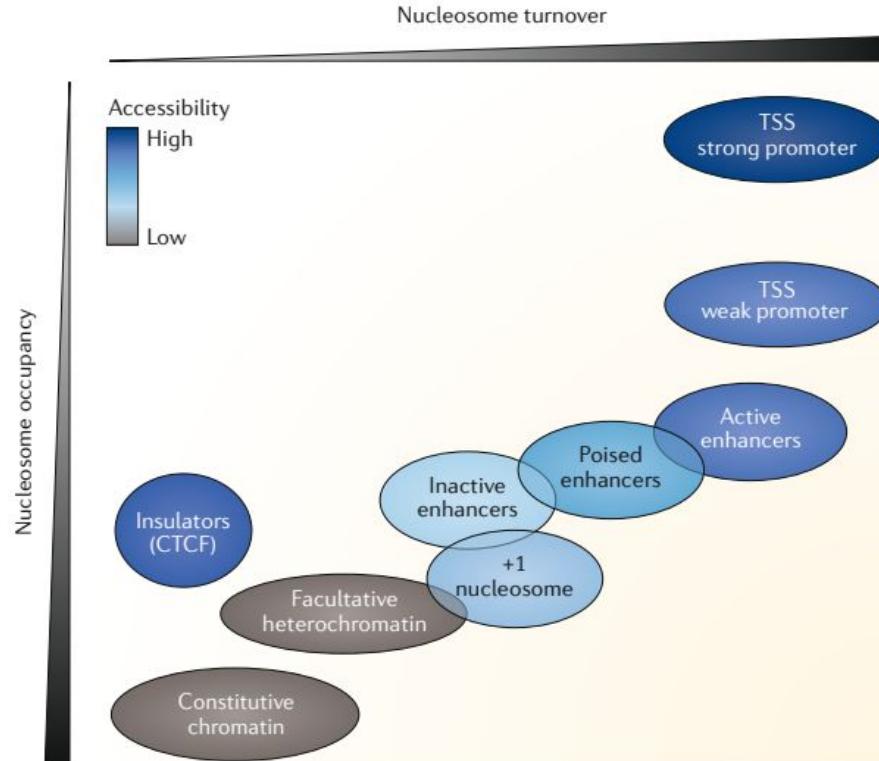
- ATAC-Seq - Assay for Transposase-Accessible Chromatin with Sequencing is a widely employed approach
- Transposase enzyme simultaneously cleaves and ligates sequencing adapters to accessible DNA



# ATAC-Seq provides maps of accessible regions of the genome (in the cells or tissues sampled)

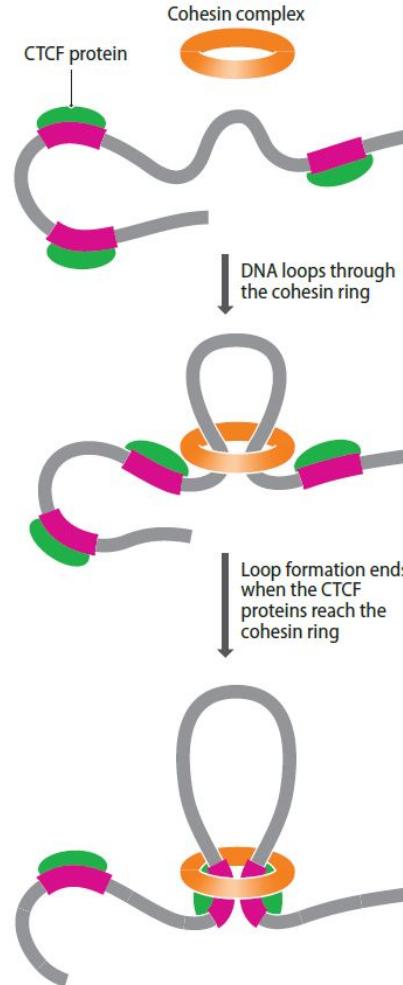


# ATAC-Seq is especially good at identifying specific elements

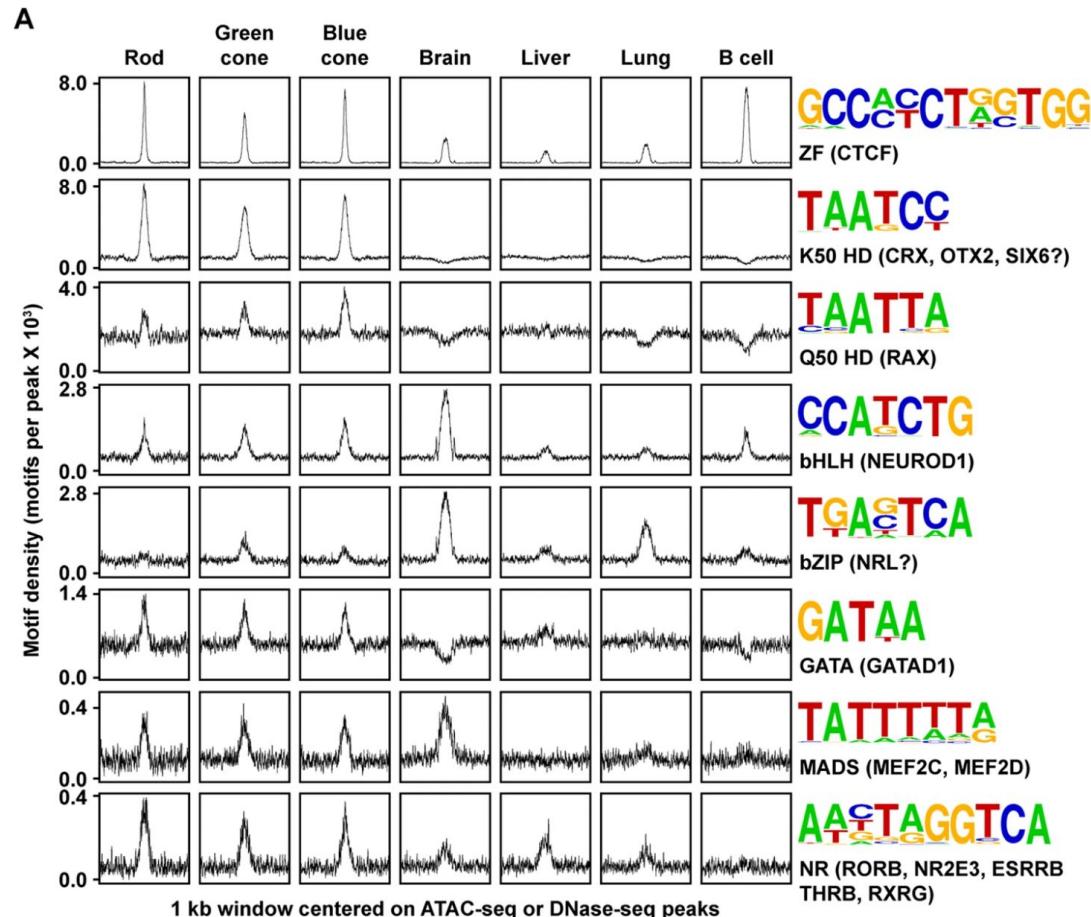


# CTCF sites define topological domains in the genome

- Regulatory interactions (activation, silencing) occur within loops.
- Such interactions generally do not occur between loops

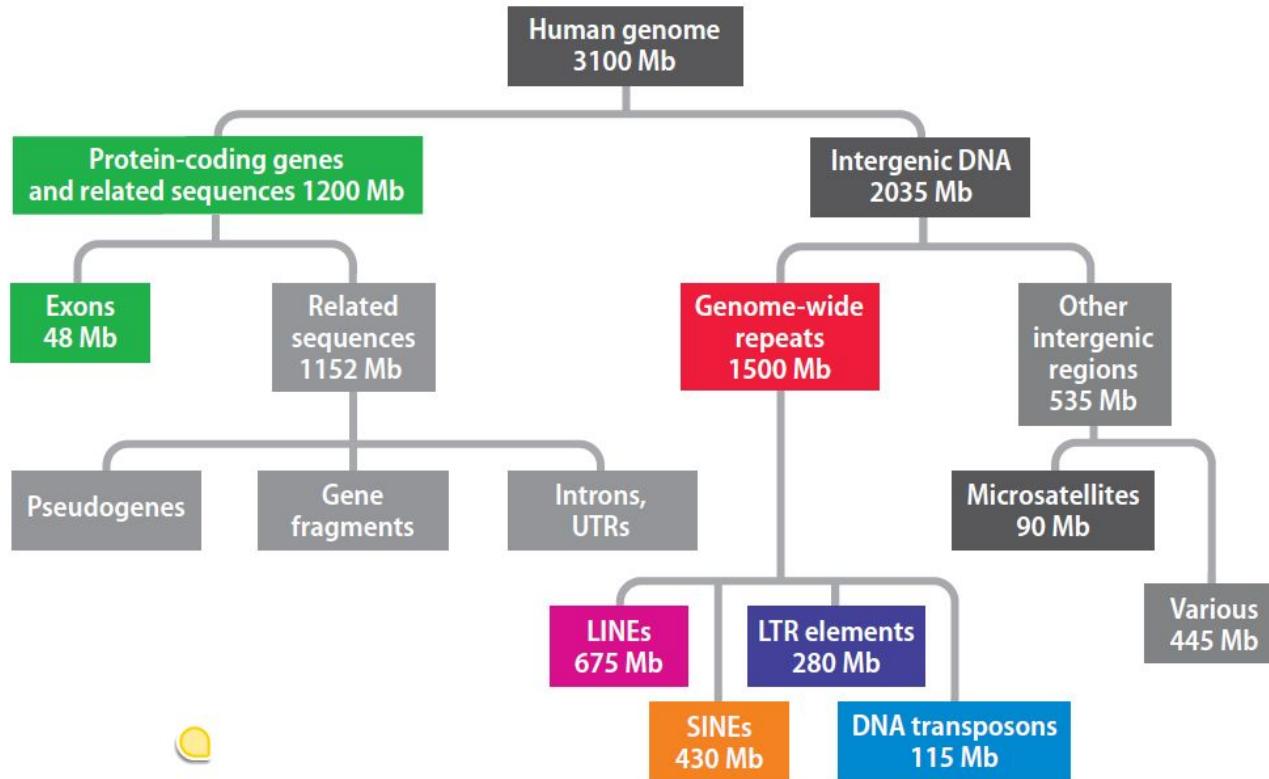


# The sequences of ATAC-seq peaks can be used to infer important transcription factors



Hughes, A. E. O., J. M. Enright, C. A. Myers, S. Q. Shen, and J. C. Corbo (2017). Cell Type-Specific Epigenomic Analysis Reveals a Uniquely Closed Chromatin Architecture in Mouse Rod Photoreceptors. *Scientific reports* 7:43184.

# Genetic features of the eukaryotic nuclear genome



# The number of protein coding genes is surprisingly small!

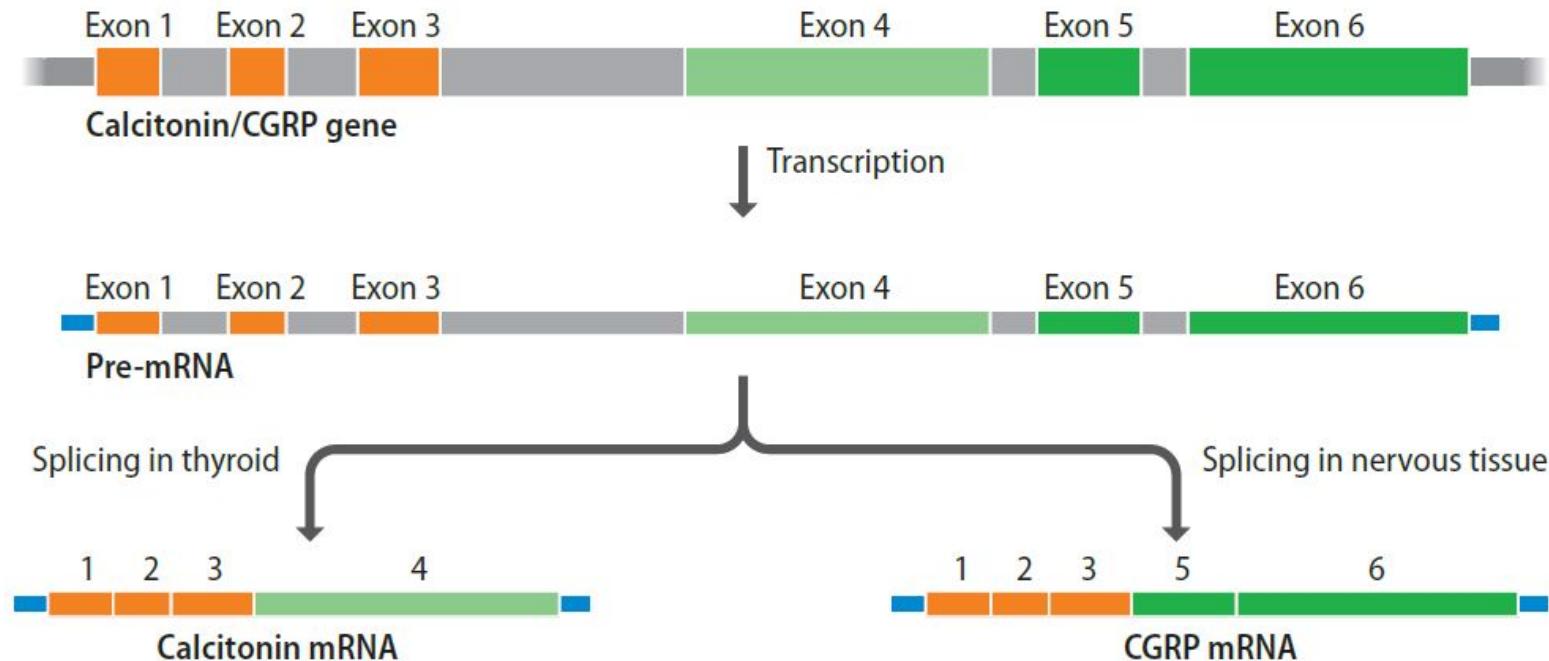
**TABLE 7.2 NUMBERS OF PROTEIN-CODING GENES FOR VARIOUS EUKARYOTES**

Species	Protein-coding genes
<i>Saccharomyces cerevisiae</i> (budding yeast)	6600
<i>Schizosaccharomyces pombe</i> (fission yeast)	5145
<i>Caenorhabditis elegans</i> (nematode worm)	20,191
<i>Arabidopsis thaliana</i> (plant)	27,655
<i>Drosophila melanogaster</i> (fruit fly)	13,968
<i>Oryza sativa</i> (rice)	37,960
<i>Gallus gallus</i> (chicken)	16,878
<i>Homo sapiens</i> (human)	20,442

Data taken from Ensembl release 104, Ensembl Plants release 51, and Ensembl Fungi release 51.

# Alternative splicing increases diversity

- In humans 20,442 genes → 78,120 proteins
- ~75% of all human protein coding genes have splice variants



# Pseudogenes

- Sequence of nucleotides that resembles a genuine gene but which does not specify a functional RNA or protein
  - Duplicated pseudogenes
    - often when gene duplication occurs one member may accumulate mutations and pseudogenize
  - Unitary pseudogene
    - single copy gene, true loss-of-function
    - e.g. *L-gulono-γ-lactone oxidase (gulo)* in Haplorrhini primates
  - Processed pseudogenes
    - reverse transcription and reintegration of mRNA sequence.
    - Most common in human genome

# Genes are unequally distributed in the genome

- In humans, ranges from:
  - 41.2 genes per Mb in Chr19
  - 6.3 in Chr13
  - just 3 per Mb in the Y chromosome
- Genes are also unevenly distributed within chromosomes:

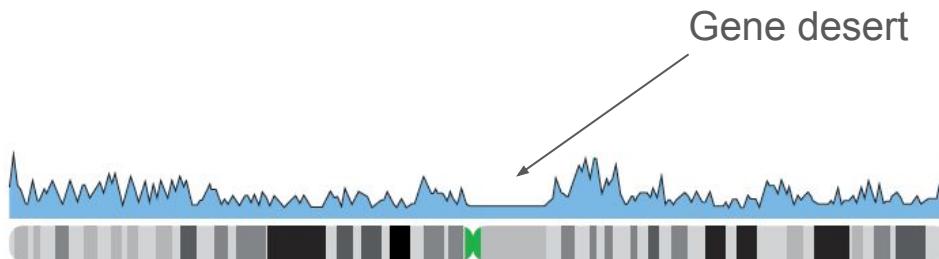
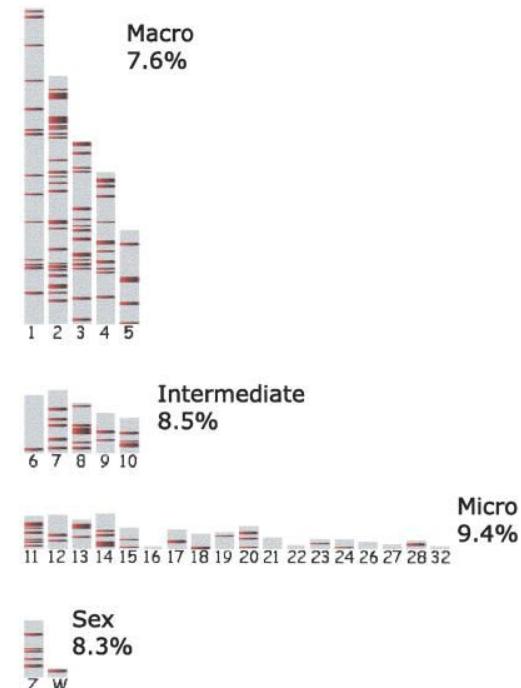


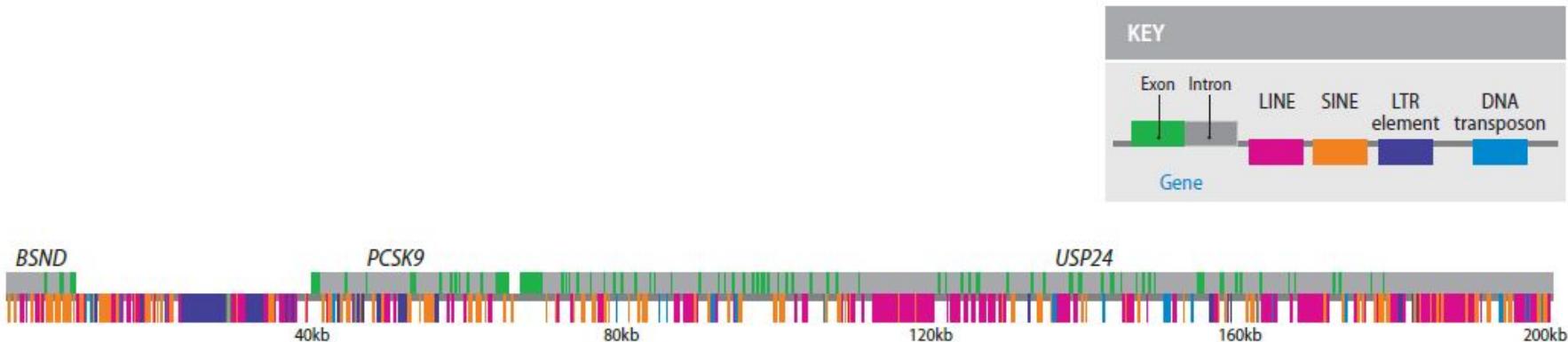
Figure 7.11 Gene density along human chromosome 1.

Gene deserts in chicken



# Protein coding genes are a very small portion of genome

- In humans, protein coding genes make up 48 Mb or ~1.5% of total sequence!
- Much of the genome (48.5%) is composed of **interspersed repeat sequences**



**Figure 7.12 A 200 kb segment of the human genome.** The segment runs from nucleotide position 55,000,000 to position 55,200,000 of chromosome 1. Within the genes, exons are shown as green boxes and introns as gray boxes. Data taken from the UCSC Genome Browser hg38 assembly.

# Interspersed repeat sequences are primarily composed of transposable elements

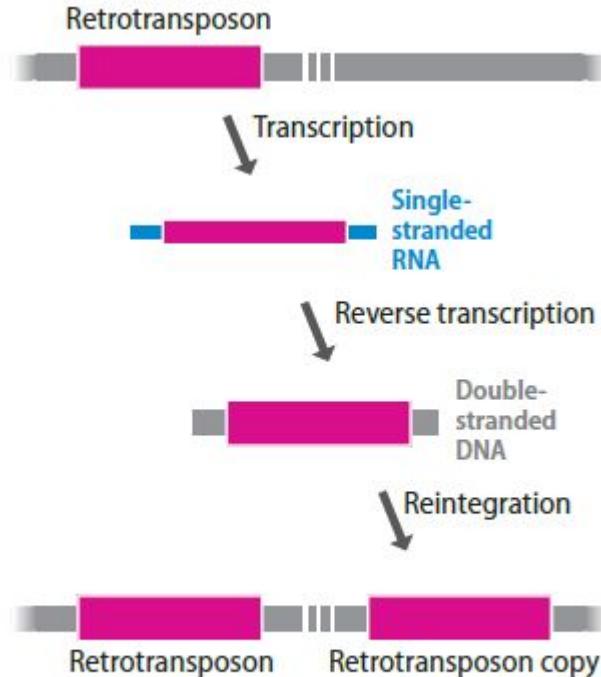
- **Transposon** - A genetic element that can move from one position to another in a DNA molecule.



**Figure 9.13** Conservative and replicative transposition.

# Interspersed repeat sequences are primarily composed of transposable elements

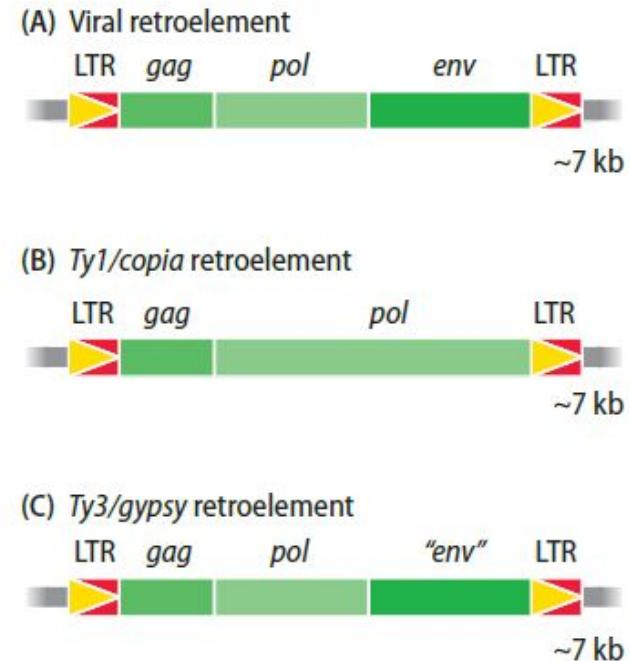
- **Retrotransposition** - transposition via an RNA intermediate



**Figure 9.14** Retrotransposition.

# Long terminal repeat (LTR) elements

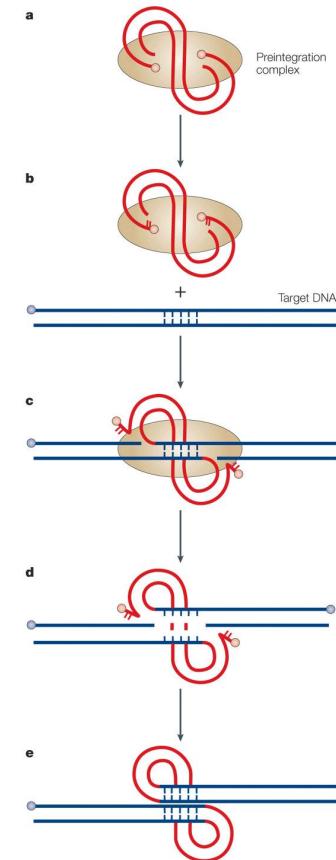
- ~8% of human genome up to 75% of maize genome
- Derive from retroviruses and encode
  - *gag* - capsid proteins
  - *pol* - reverse transcriptase
- May produce virus-like particles
- Many are “decayed” with loss of function mutations in genes



**Figure 9.15** Genome structures for LTR

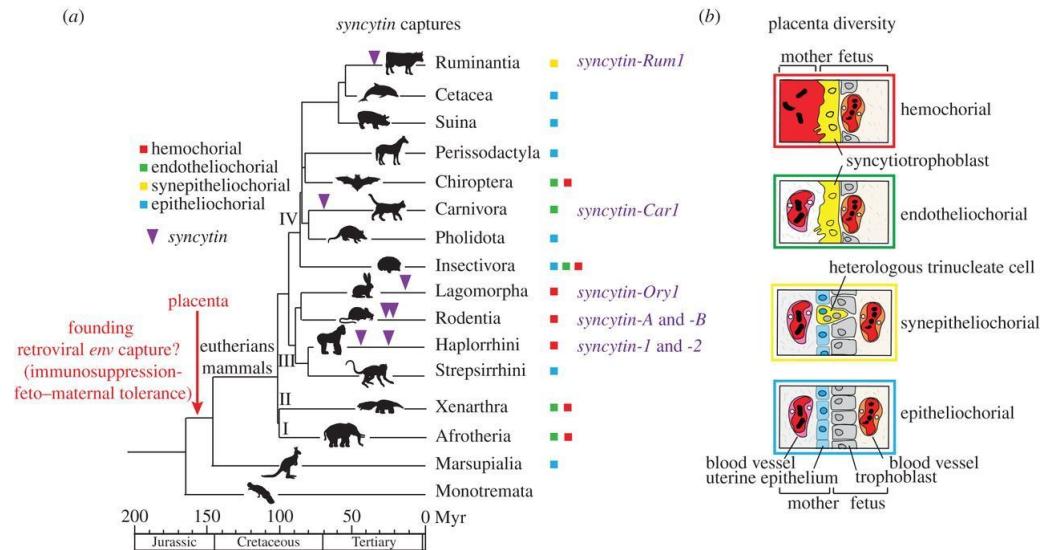
# Long terminal repeat (LTR) elements

- ~8% of human genome up to 75% of maize genome
- Derive from retroviruses and encode
  - *gag* - capsid proteins
  - *pol* - reverse transcriptase
- May produce virus-like particles
- Many are “decayed” with loss of function mutations in genes
- LTR sequences mediate integration



# Syncytins are derived from retroviruses

- Syncytins are proteins essential for placental development and implantation in mammals
- Evolved from endogenous retroviral envelope (*env*) gene



# Short interspersed nuclear elements (SINEs)

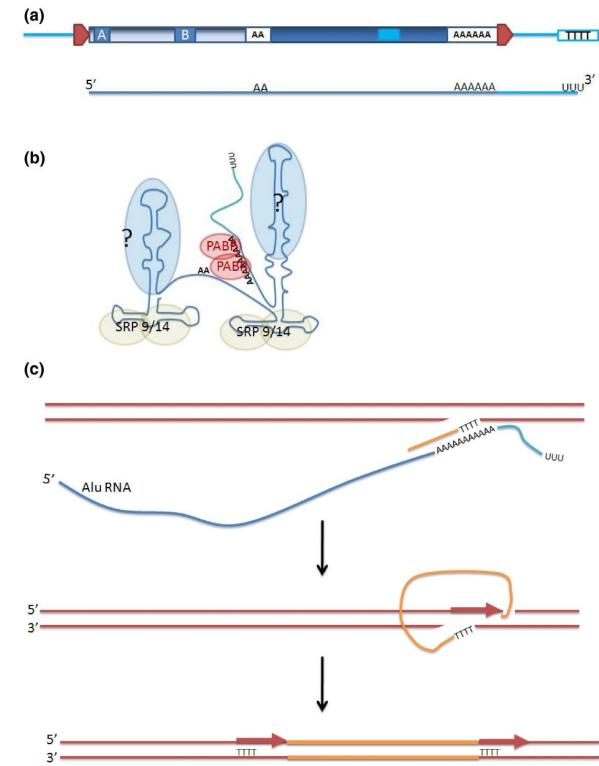
- 1.7 million copies in human genome (~14%)
- Nearly 10% of human genome *Alu* elements
  - ~280 bp
  - Derives from 7SL RNA gene
  - Requires other elements (LINEs) for expression and replication.
  - Provide homology across genome and may drive high levels of recombination
  - Varied effect on gene expression and function

## SINEs (~ 100-300 bp)

### 7SL RNA head



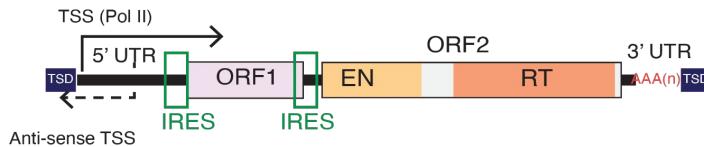
Alu - 280bp (human)



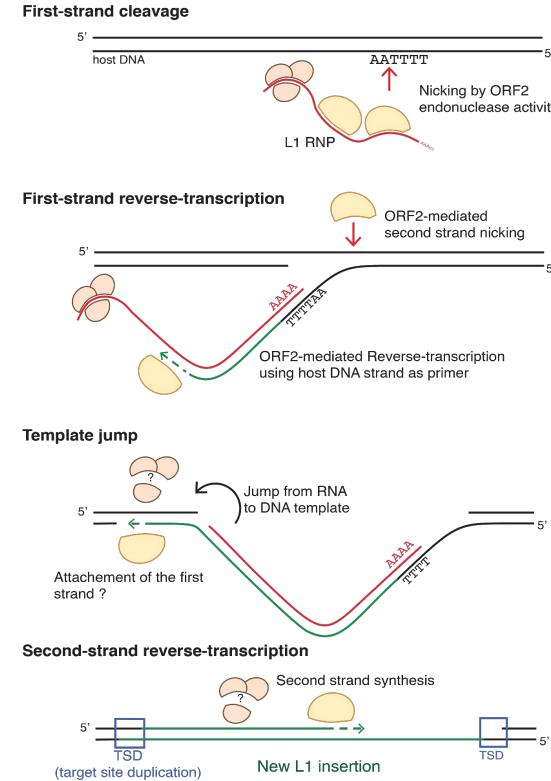
# Long interspersed nuclear elements (LINEs)

- ~6 kb in length, but many truncated
- >20% of human genome
- *gag* and *pol* elements, but not LTRs

Mouse LINE1 (~ 6 kb)

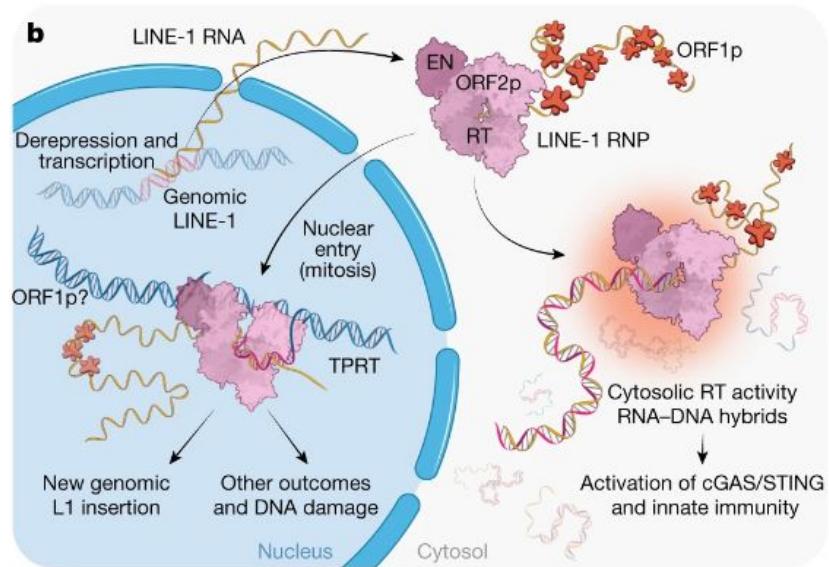


**EN:** Endonuclease domain. **RT:** Reverse transcriptase domain.  
**IRES:** Internal ribosomal entry site **TSD:** Tandem site duplication.  
**A** and **B:** RNA pol III promoter box



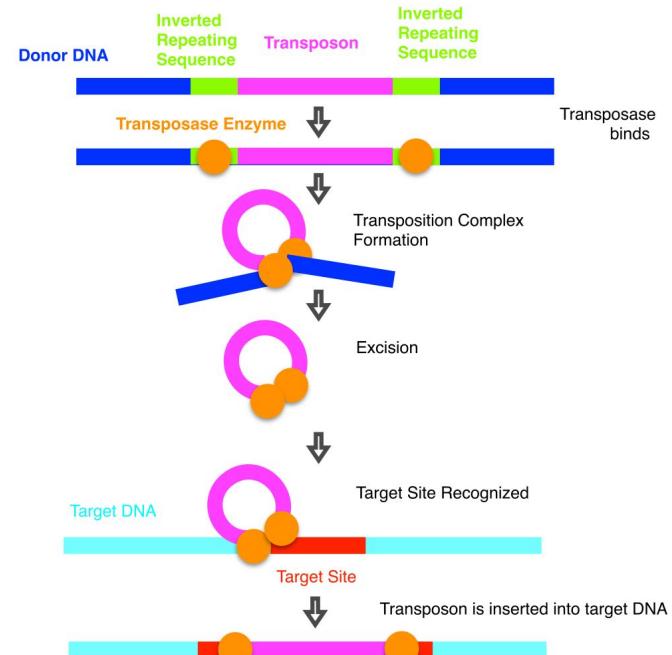
# Long interspersed nuclear elements (LINEs)

- ~6 kb in length, but many truncated
- >20% of human genome
- *gag* and *pol* elements, but not LTRs
- Activation of LINE-1 expression implicated in disease processes



# DNA transposons

- 3.7% of human genome, most in active
- More abundant in plants
- First described by Barbara McClintock in 1944



[https://en.wikipedia.org/wiki/DNA\\_transposon#/media/File:Cut\\_and\\_Paste\\_mechanism\\_of\\_transposition.svg](https://en.wikipedia.org/wiki/DNA_transposon#/media/File:Cut_and_Paste_mechanism_of_transposition.svg)

# Microsatellites

- ~3% of human genome
- Very short (<13 bp) tandemly repeated sequences
- Produced through polymerase slippage during DNA replication
- Highly variable - lengths are different among all individuals
- Major mechanism of “DNA fingerprinting”

