

GENOMA HUMANO

Prof. Leandro Martins de Freitas, PhD
IMS/UFBA



HISTÓRICO



Complexidade dos
organismos é muito grande



Complexidade da formação
do fenótipo é muito grande



Muitos modelos
hereditários e genéticos

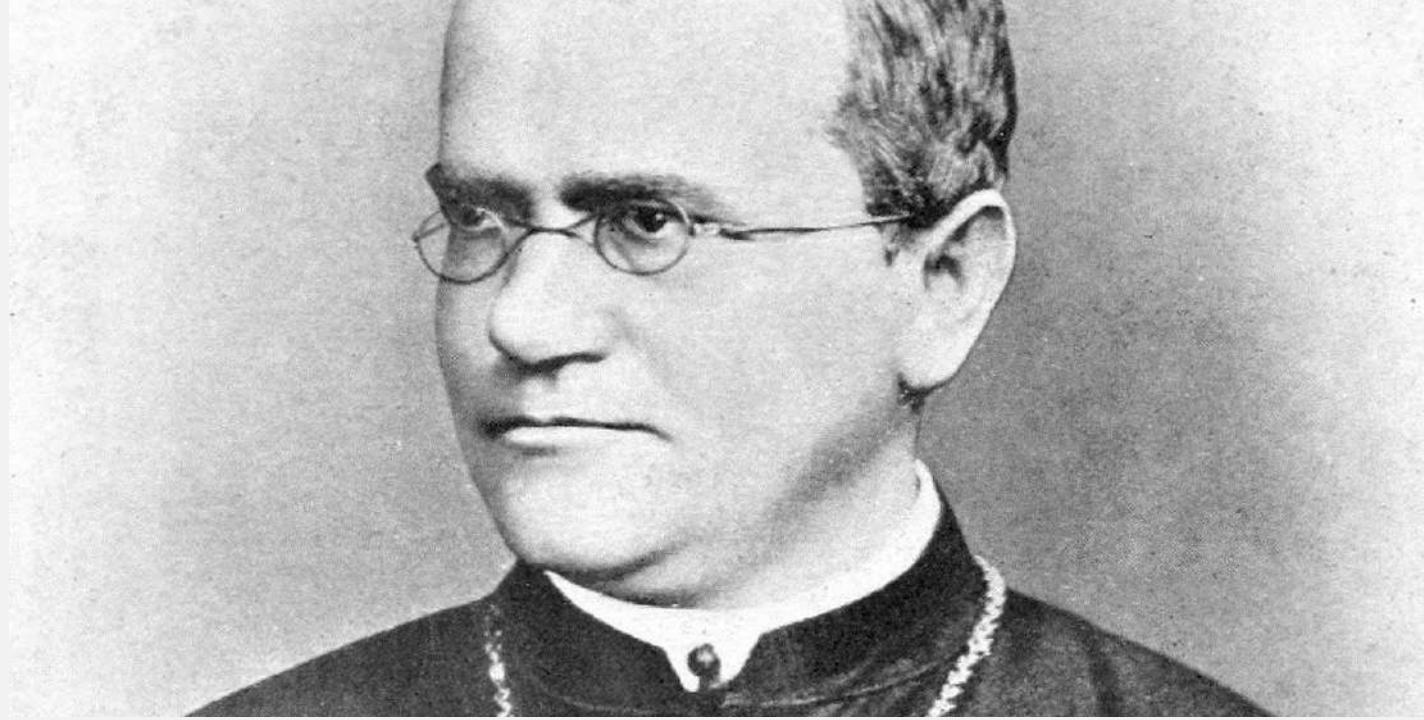
HISTÓRICO

Genes não atuam sozinhos -
Interação dos genes

Conhecer somente a sequência
de um gene não é suficiente

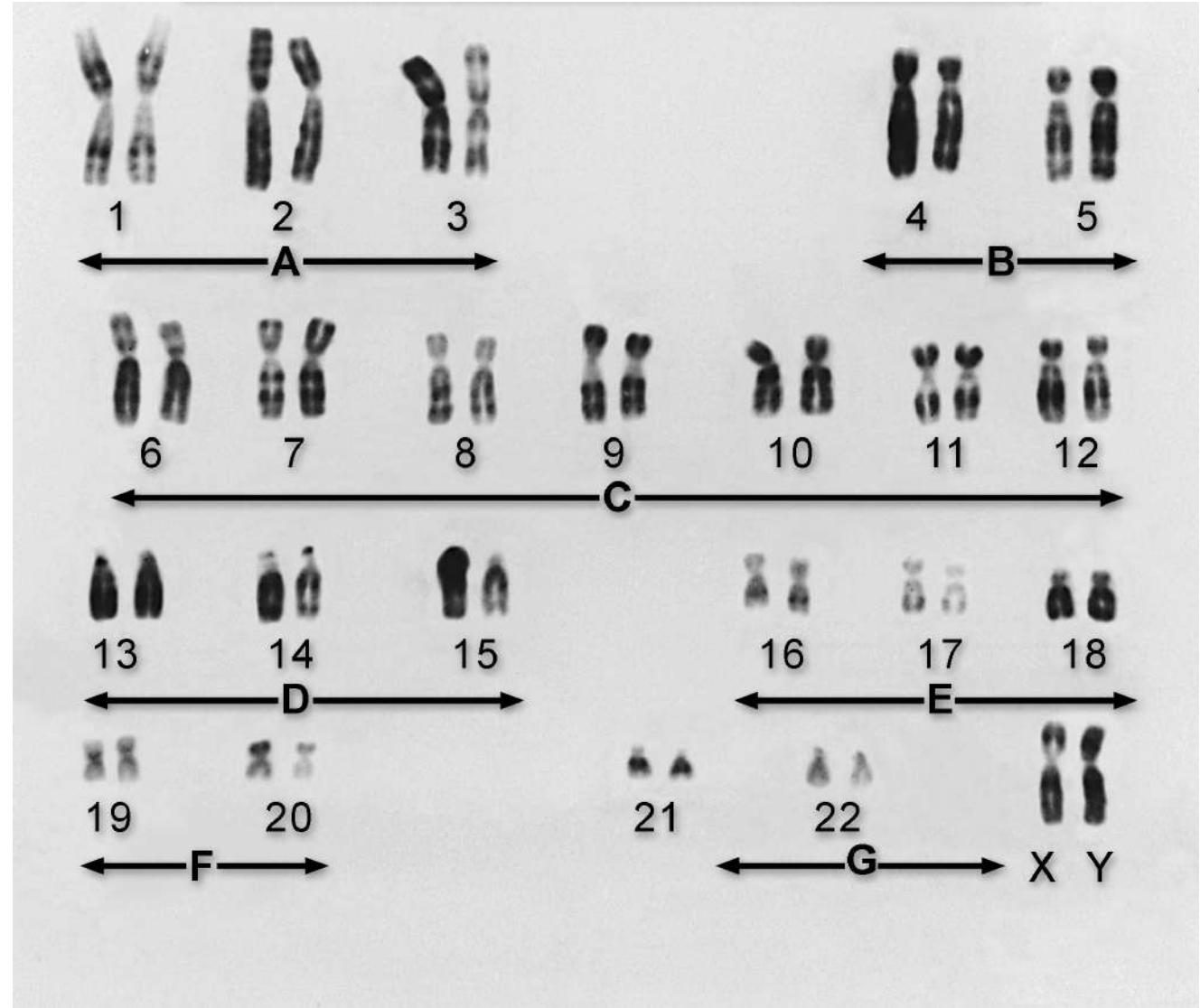
Necessidade de conhecimento
da informação genética

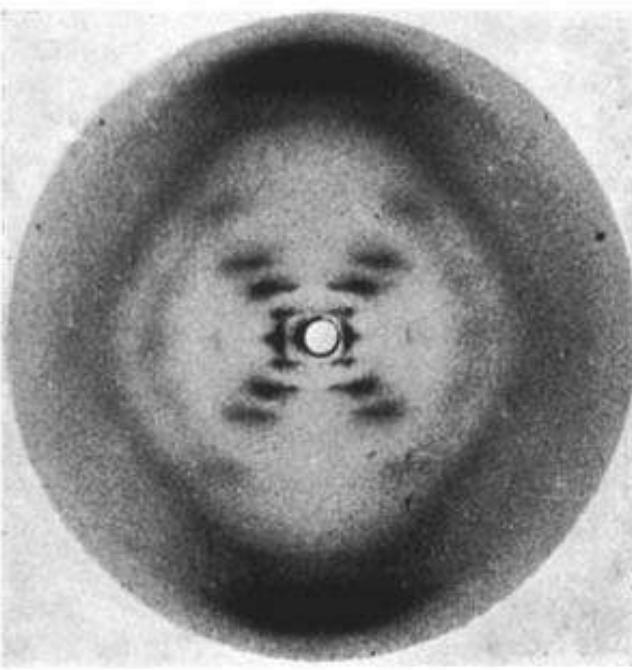
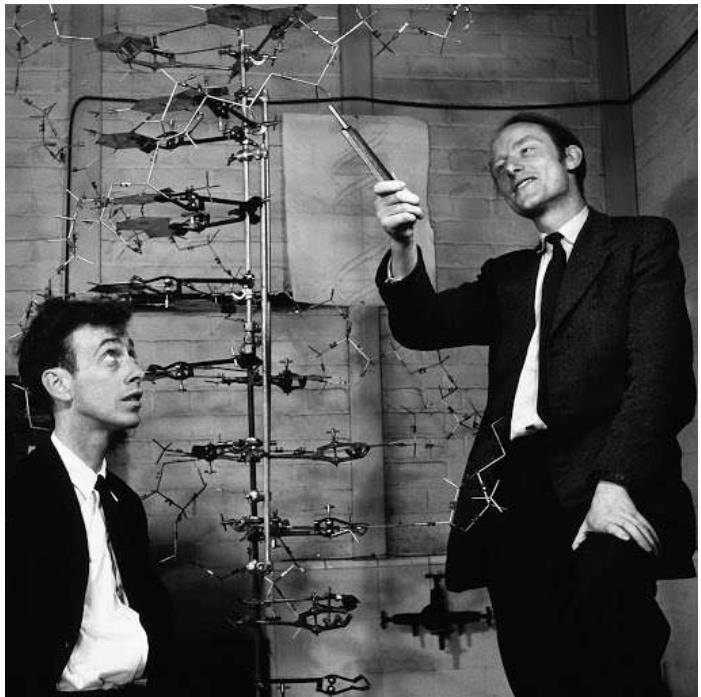
GENÉTICA



Characteristics of pea plants Gregor Mendel used in his inheritance experiments						
Seeds	form		Flower colour	Pod		Stem position of inflorescences
	cotyledons	form		colour	size	
round roundish			yellow			axial
	wrinkled			violett-red	yellow	
wrinkled			white			long
			green	constricted between the seeds	green	
terminal						short
						

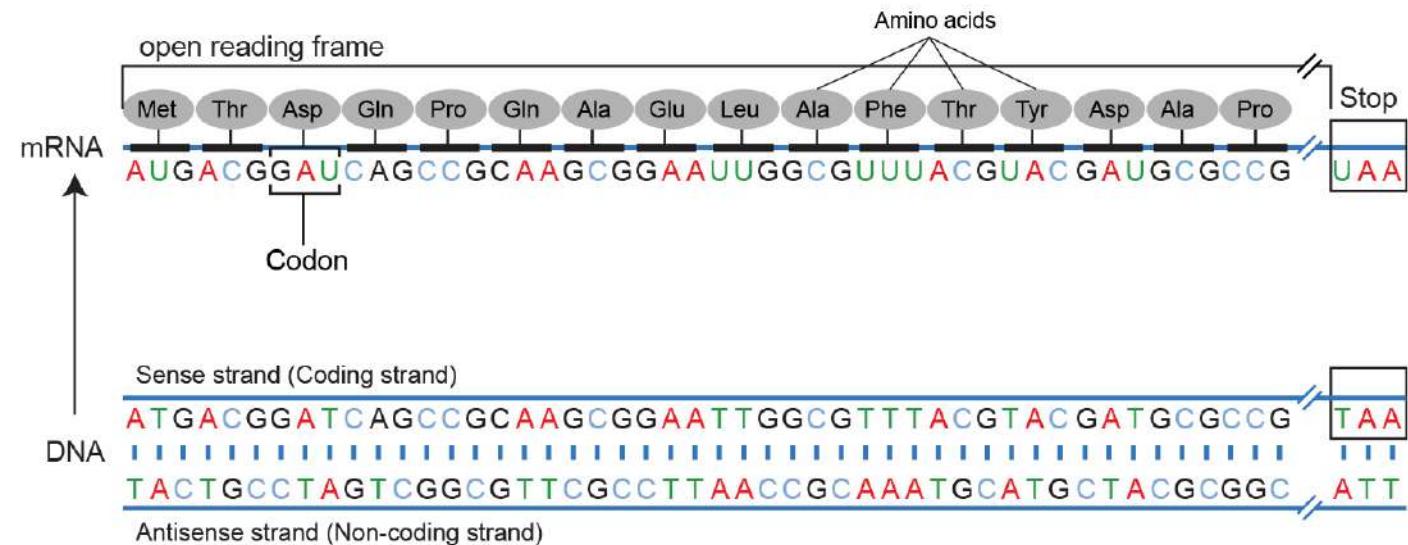
CROMOSSOMOS



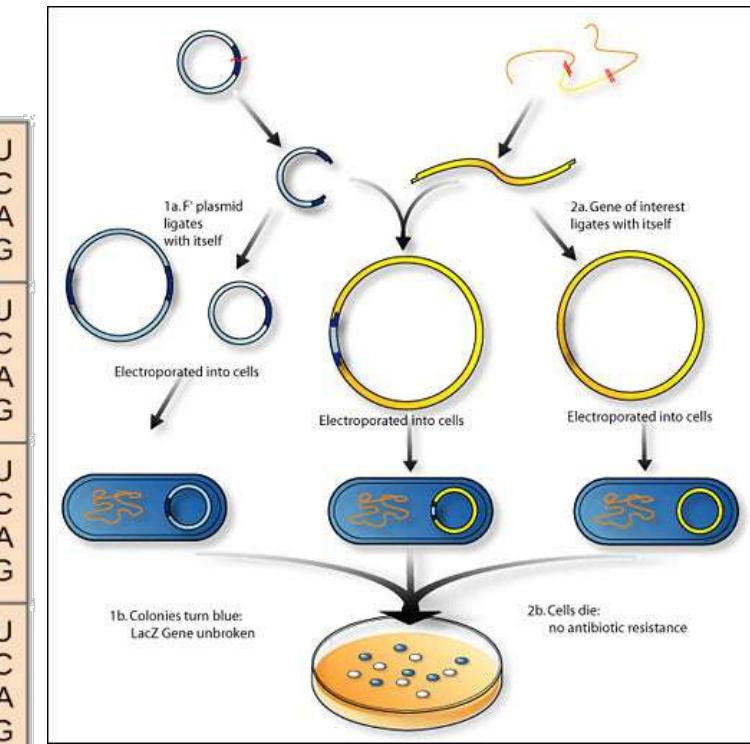


DNA – DUPLA HÉLICE

LEITURA DO DNA



		Second letter				
		U	C	A	G	
First letter	U	UUU } Phe UUC } UUA } Leu UUG }	UCU UCC UCA UCG }	UAU } Tyr UAC }	UGU } Cys UGC }	UCA A G
	C	CUU } CUC } Leu CUA } CUG }	CCU CCC CCA CCG }	CAU } His CAC }	CGU } Arg CGC CGA CGG }	U C A G
	A	AUU } Ile AUC } AUA }	ACU ACC ACA ACG }	AAU } Asn AAC }	AGU } Ser AGC }	U C A G
	G	AUG } Met	GUU GUC GUA GUG }	GCU GCC GCA GCG }	GAU } Asp GAC }	GGU } Gly GGC GGA GGG }



GENÔMICA

A word cloud graphic on a white background containing various genomic and bioinformatics terms. The words are colored in shades of red, green, yellow, and orange, and are arranged in a non-linear, overlapping pattern. Some words have smaller descriptive text next to them.

Variant
samples
Bioinformatics
Microarrays
Genotyping
Genomics
Analysis
Sequencing
NGS
miRNA
Exome
nCounter
NanoString
Tiling
Real-time
discrimination

CNV
Illumina
Caliper
allelic
profiling
ChIP-seq
expression analysis
Qubit
Transcriptomics
Affymetrix
Q-PCR
TaqMan
ChIP-string
BeadStation
Sybr
Epigenetics
RNA-seq
Gene
Methylation
RNA
Gene
fluorescence-luminescence
Covaris
PCR
cBot
SNP
arrays
Green
sequencing
FFPE
Bioanalyzer
arrays
DNA
TapeStation
Real-time
discrimination

PRIMEIROS GENOMAS

Virus bacteriano Φ X174 – 1977,
1978

Virus bacteriano Lambda – 1982

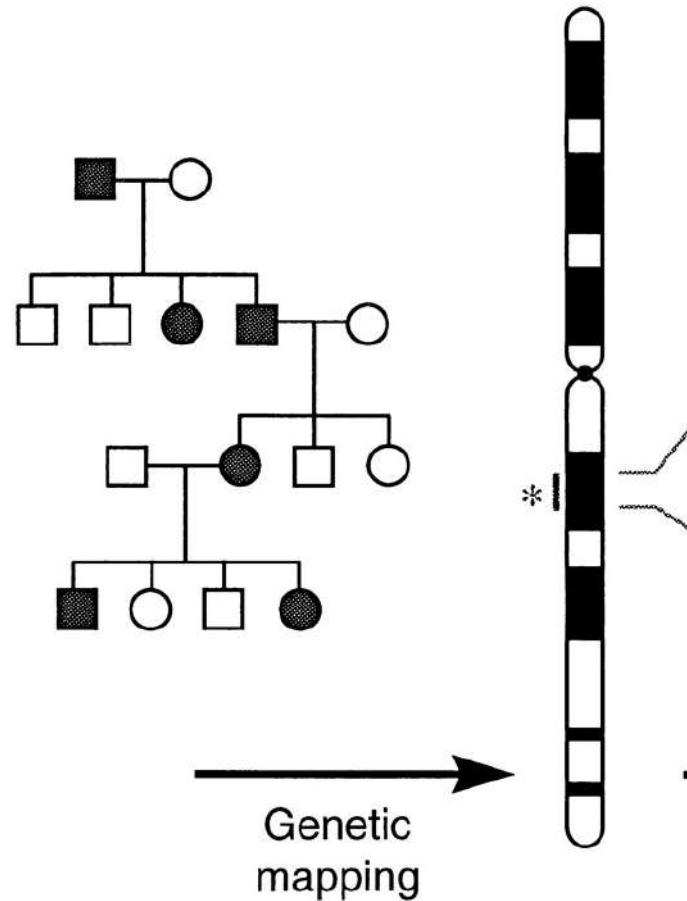
Vírus SV40 -1978

Mitocôndria humana 1981

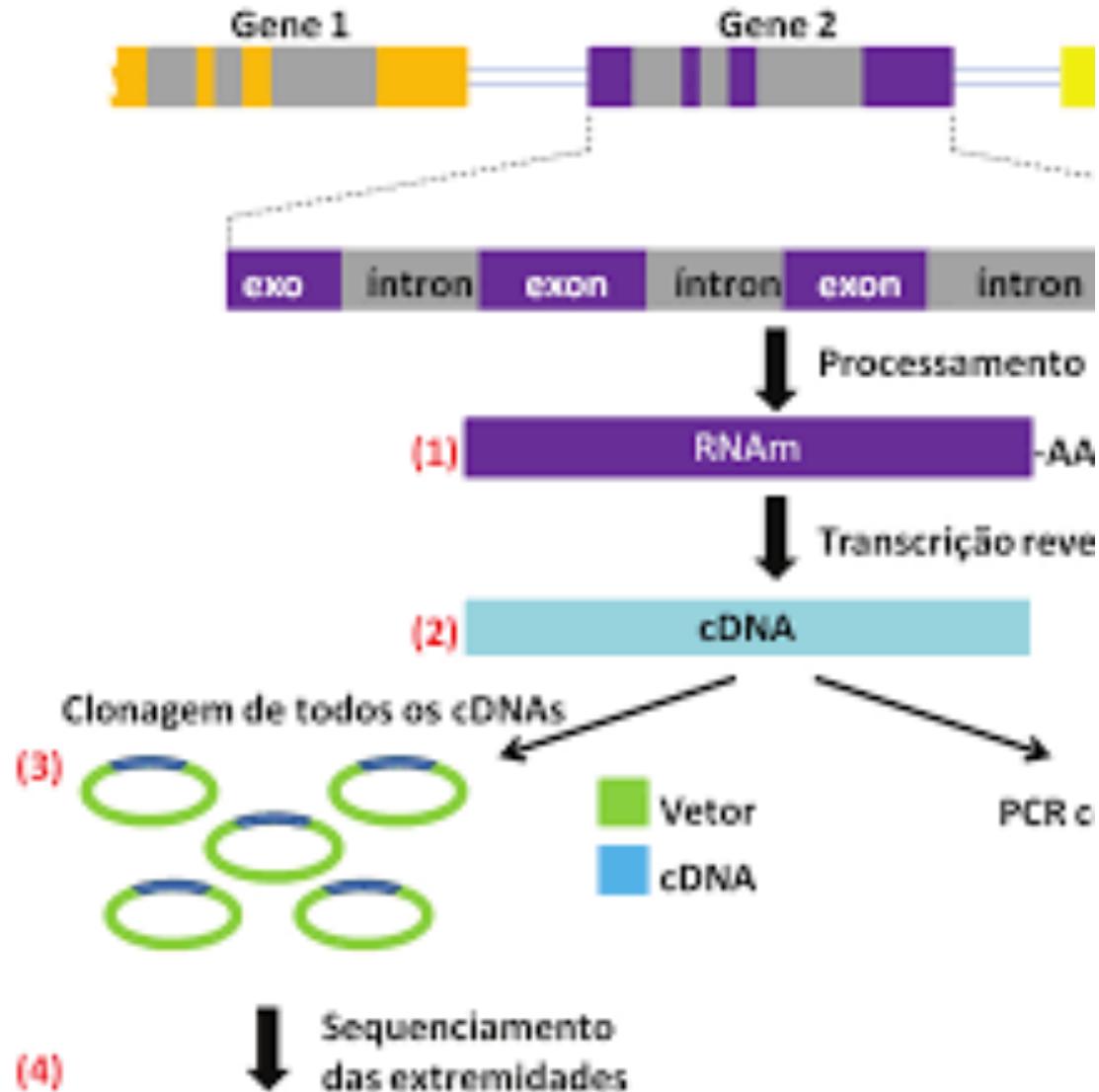
MAPA GENÉTICO

Family
studies

Chromosome
interval



SEQUENCIAMENTO SHOTGUN DE cDNA



GENOMA HUMANO

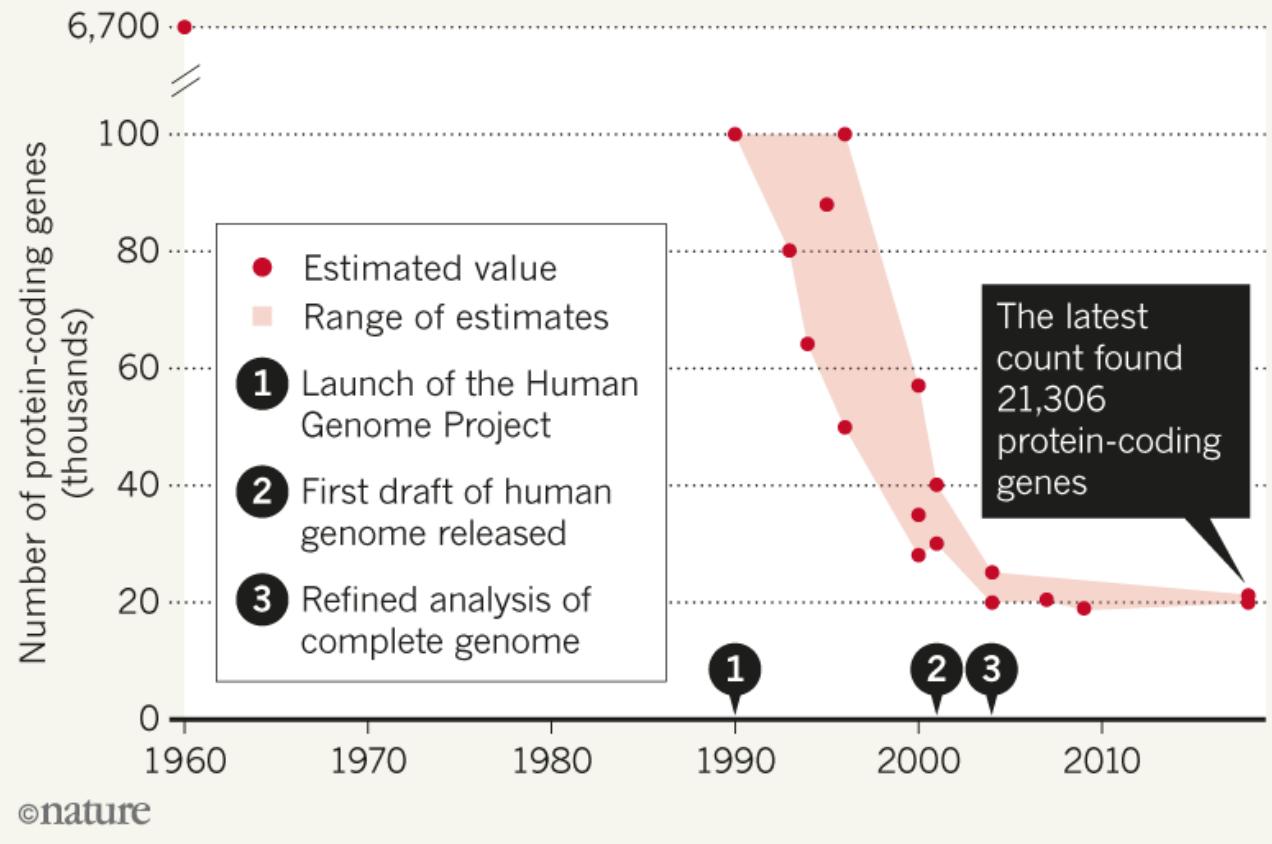
30-40 mil genes, estimativas no início

96% das regiões de eucromatina foi sequenciado (94% do genoma)

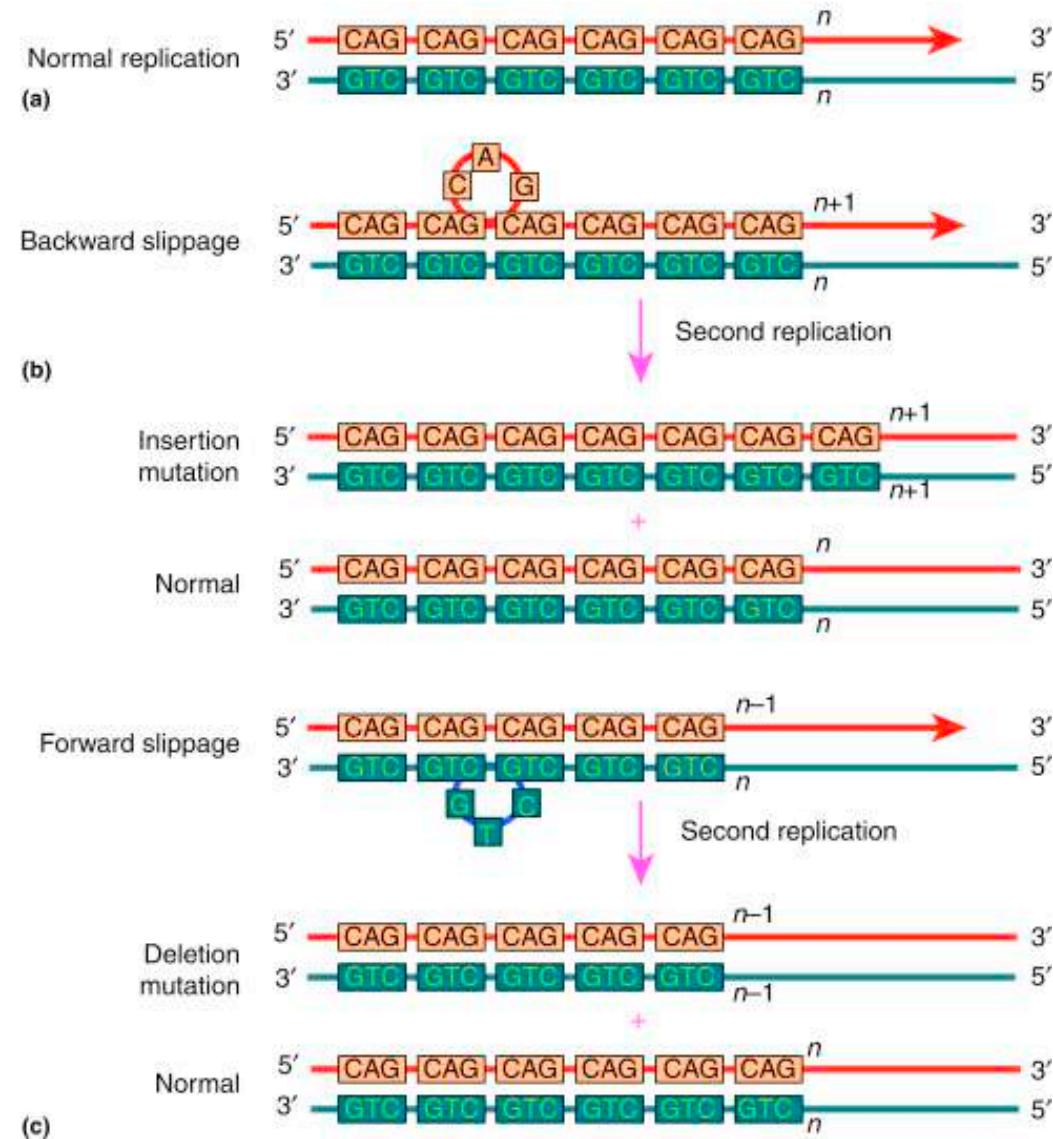
GENES HUMANOS - PROTEINAS

GENE TALLY

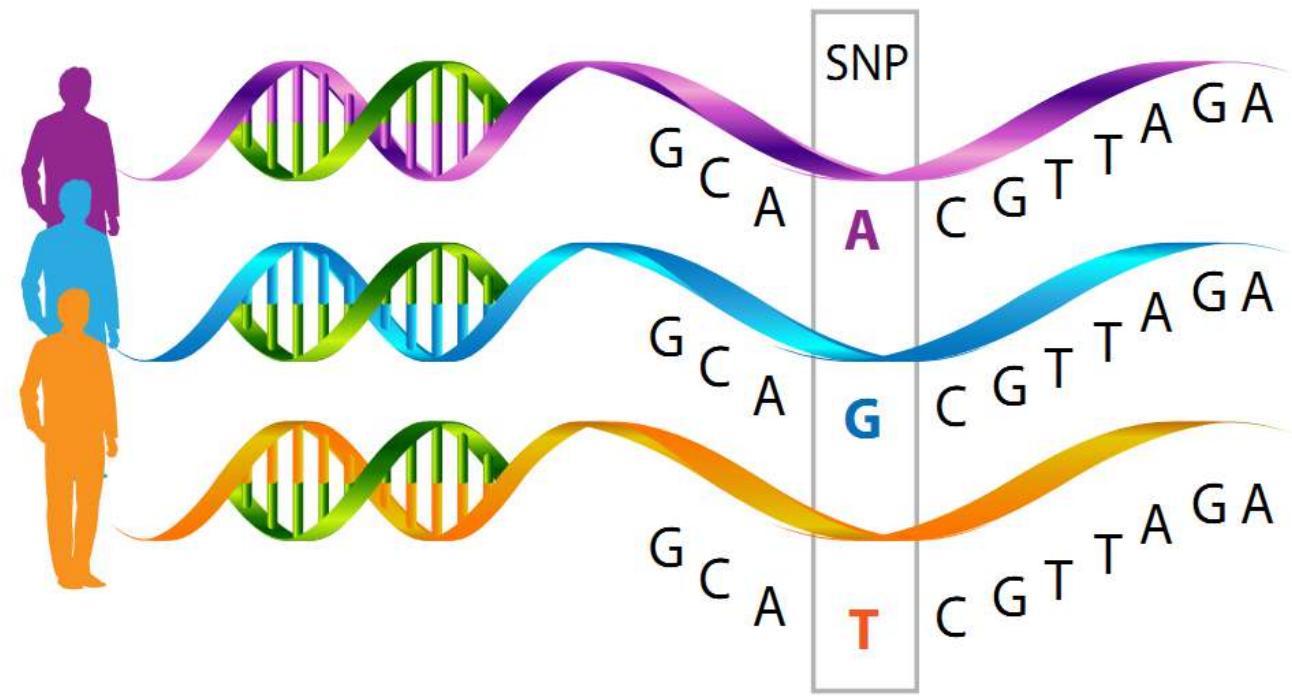
Scientists still don't agree on how many protein-making genes the human genome holds, but the range of their estimates has narrowed in recent years.



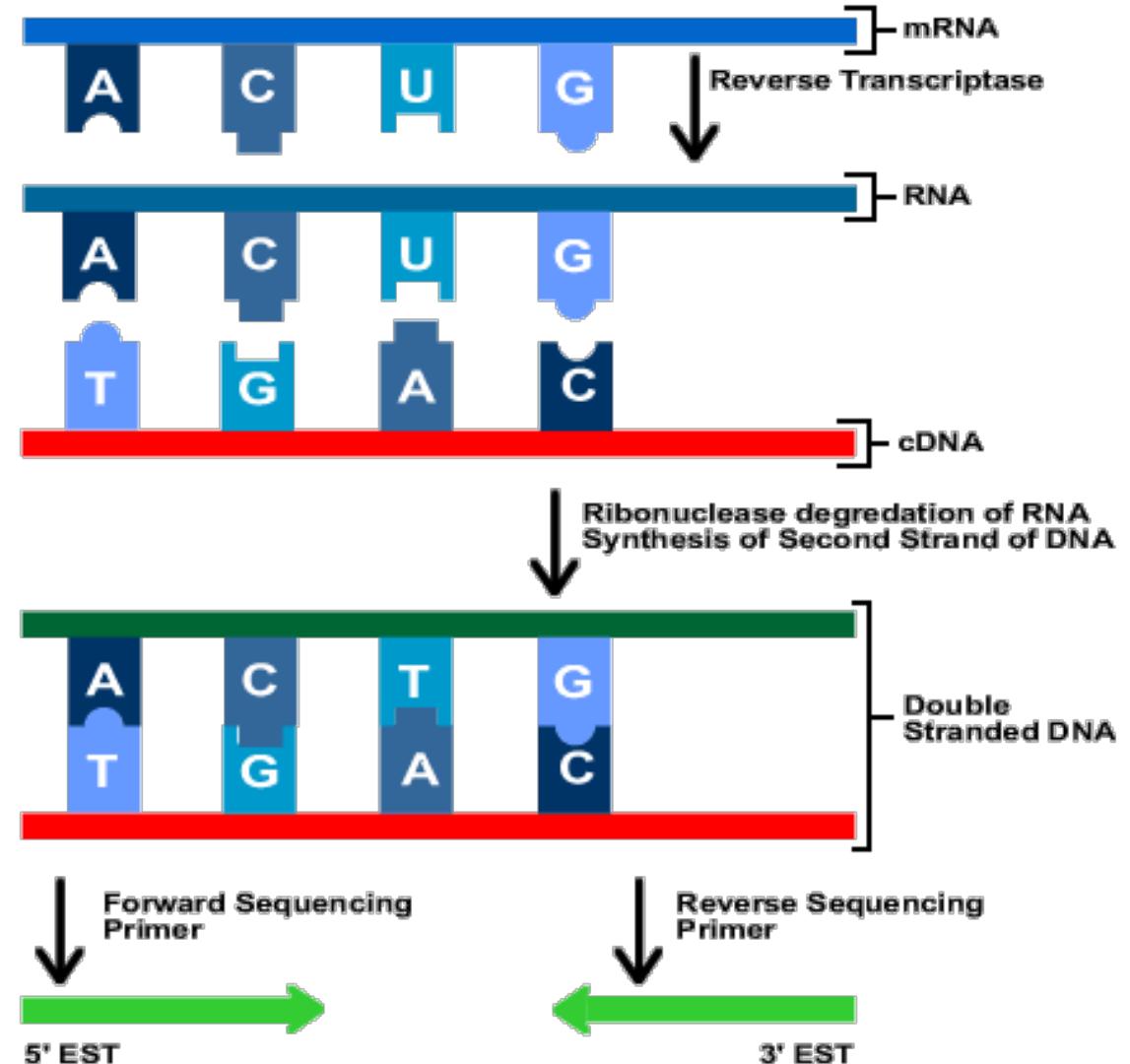
MICROSSATÉLITES

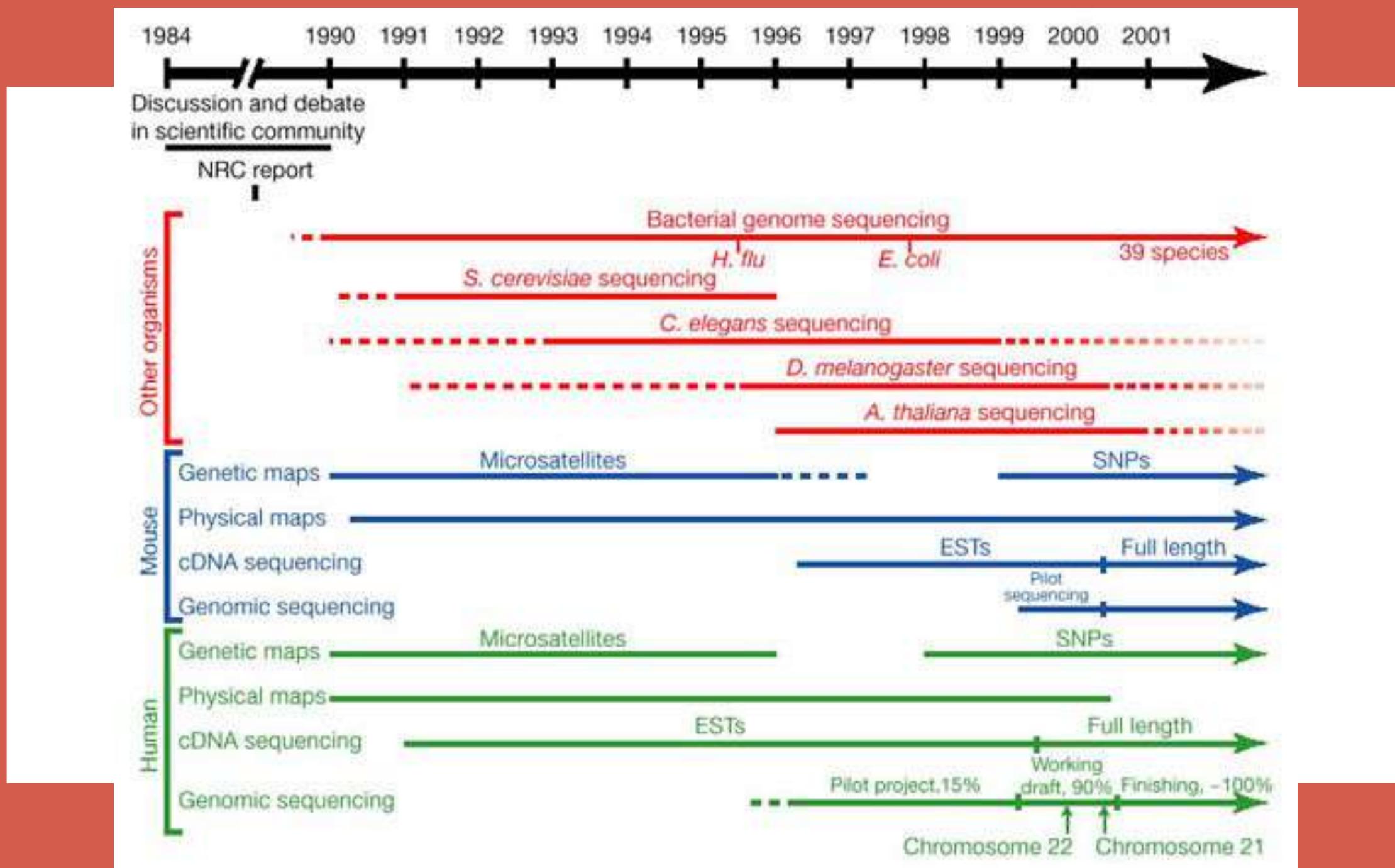


SNP

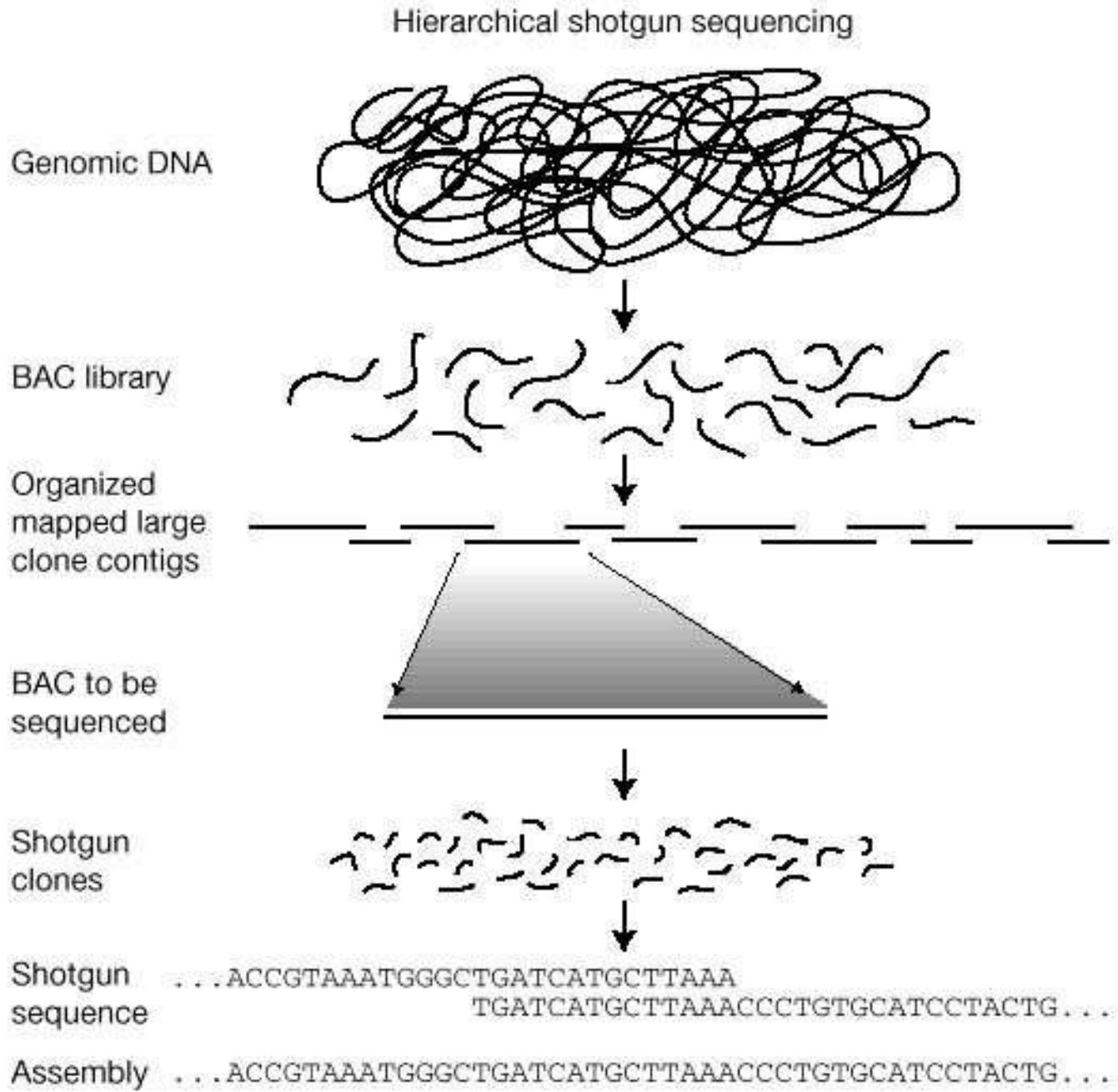


EST - ESPRESSION SEQUENCE TAG



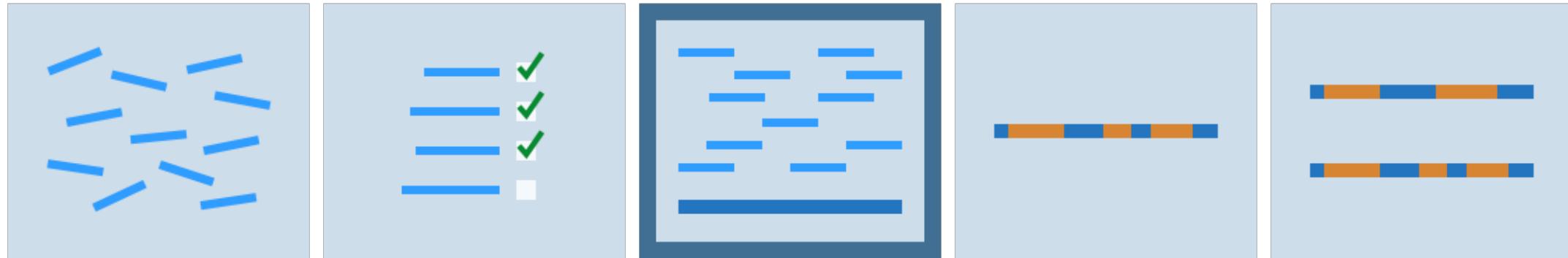


GENOMA - SHOTGUN



PRODUÇÃO AUTOMATIZADA - CENTER FOR GENOME RESEARCH





Sequencing

Quality control

Assembly

Annotation

Comparison

GENOMA E BIOINFORMÁTICA

ESTRATÉGIAS DE MONTAGEM

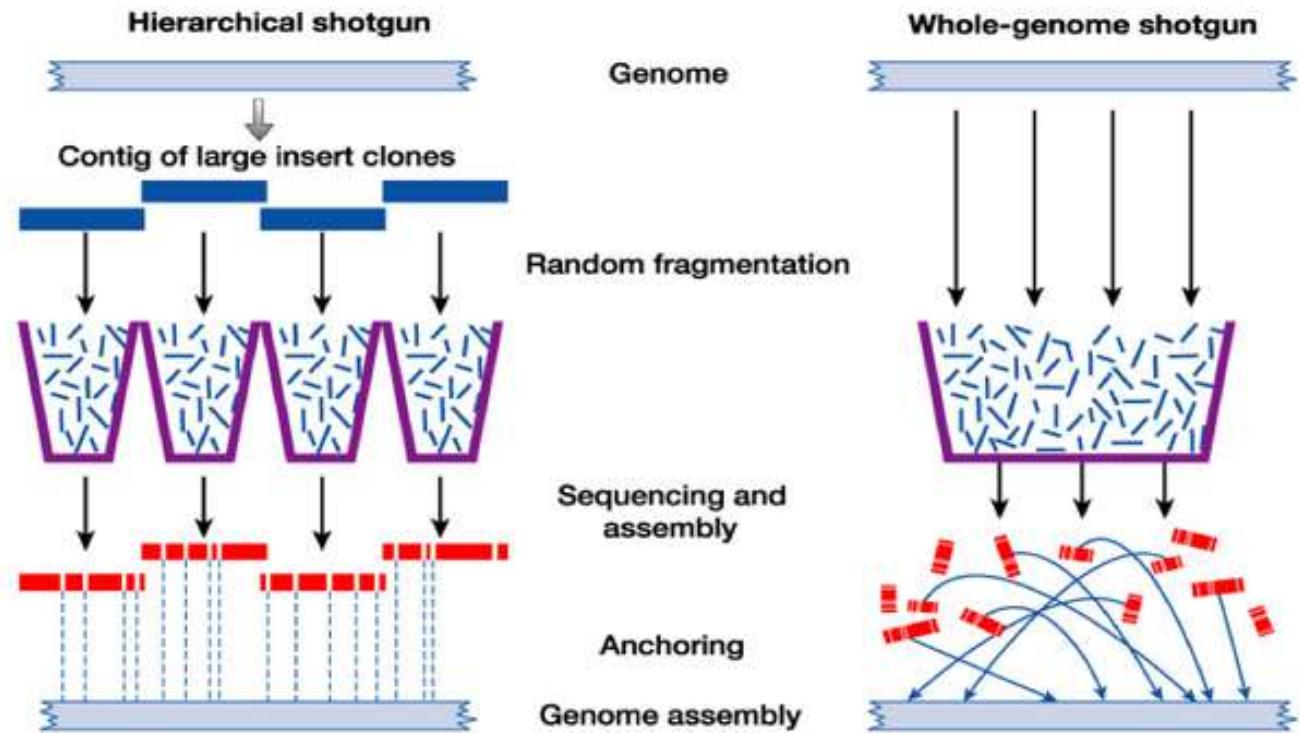


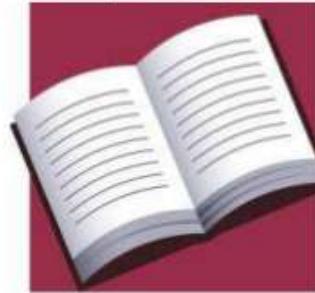
Figure 8-3 Human Molecular Genetics, 3/e. (© Garland Science 2004)

ANALOGIA DO GENOMA

A single recipe is like...



A recipe book is like...



Two copies of 23 recipe books is like...



A Gene:

One set of instructions for how to make one protein.

A Chromosome:

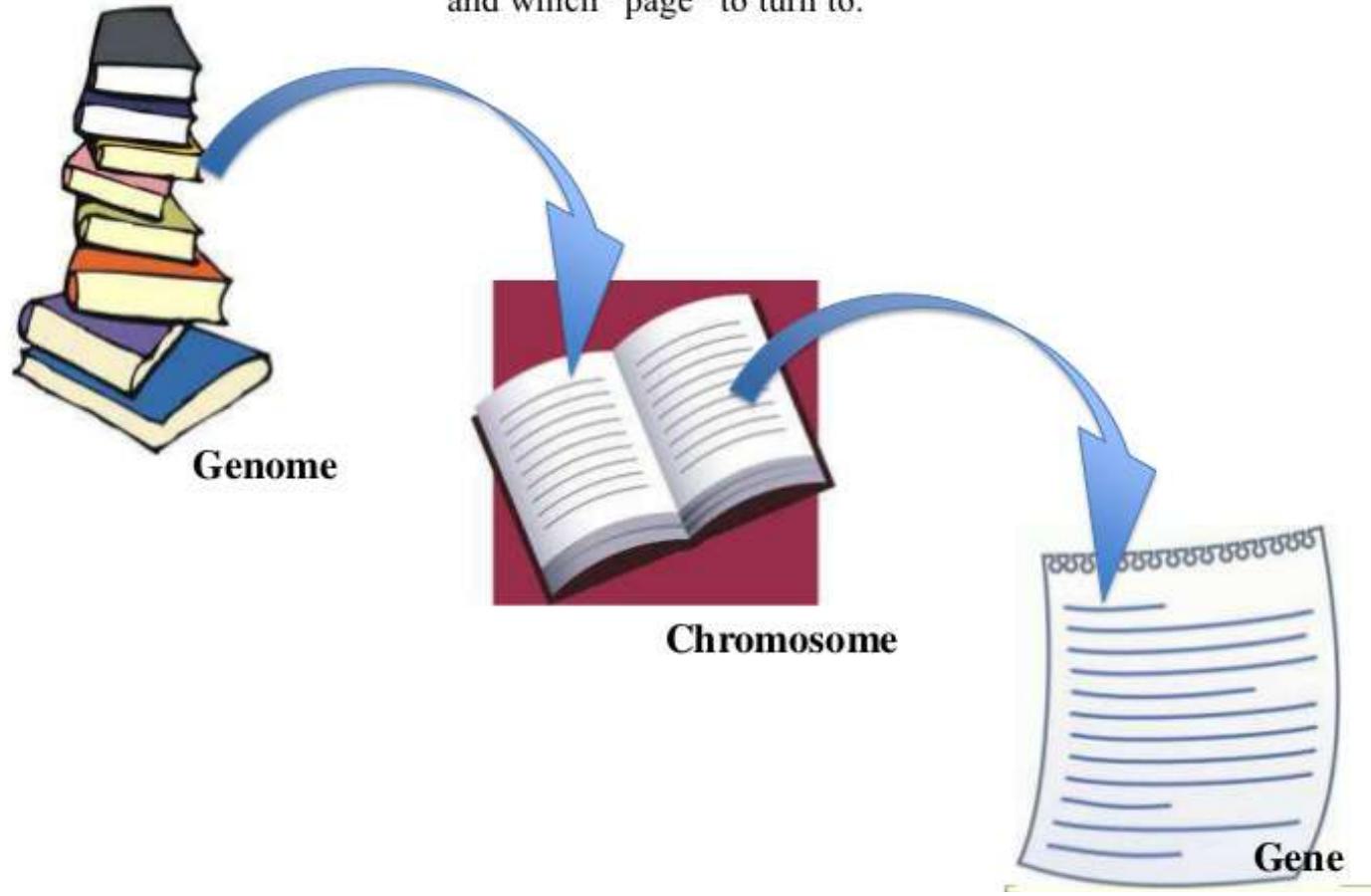
Thousands of sets of instructions for how to make thousands of proteins

Our Genome:

ALL of the sets of instructions for how to make ALL the proteins we need

ALL (gene, chromosome, genome) are written in same the DNA alphabet!!!

In order for scientists to find a specific gene, they need to know which “book” to look in,
and which “page” to turn to.



MONTAGEM DE UM GIGANTESO QUEBRA- CABEÇA



GIGANTESCO QUEBRA- CABEÇA SEM A CAIXA

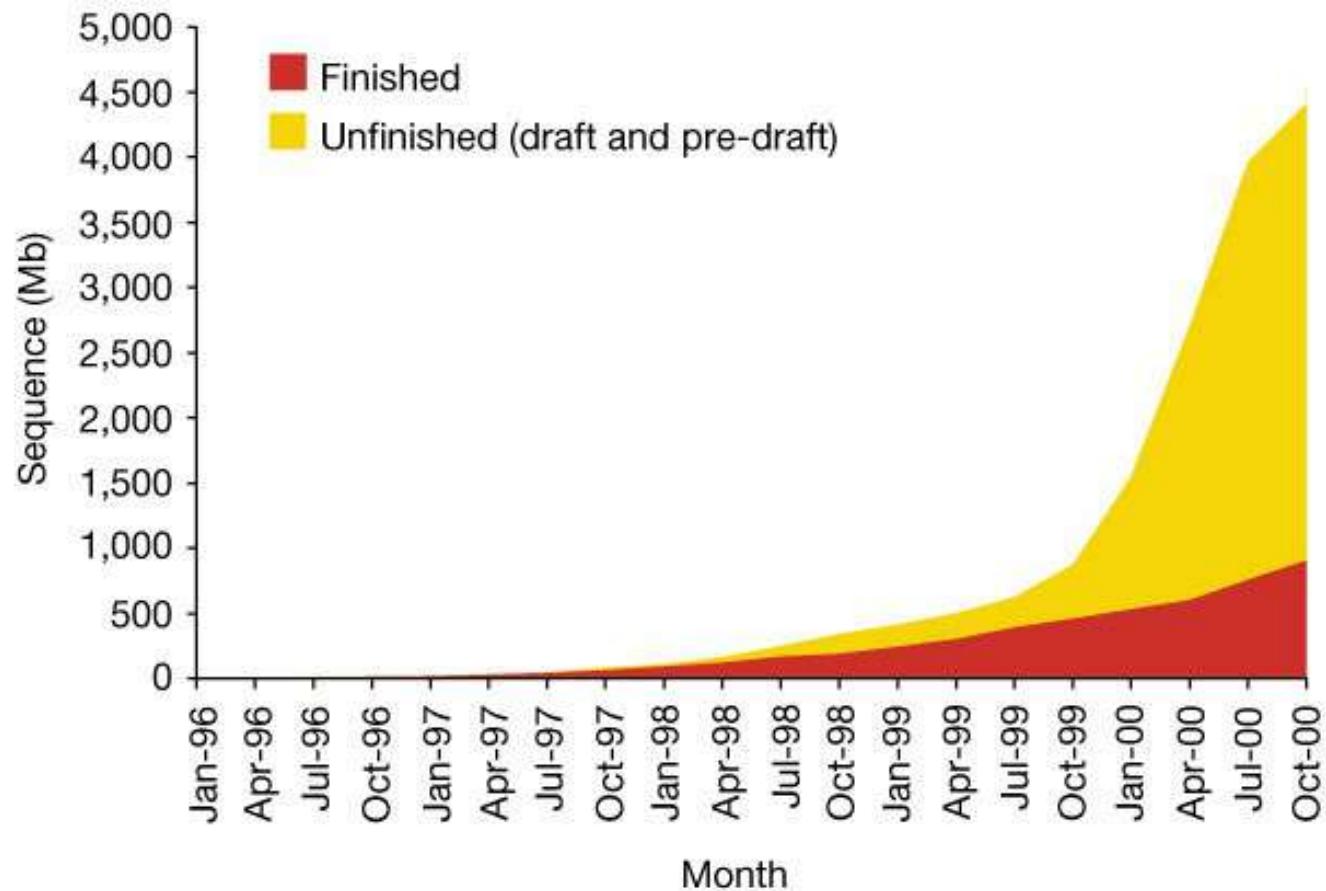




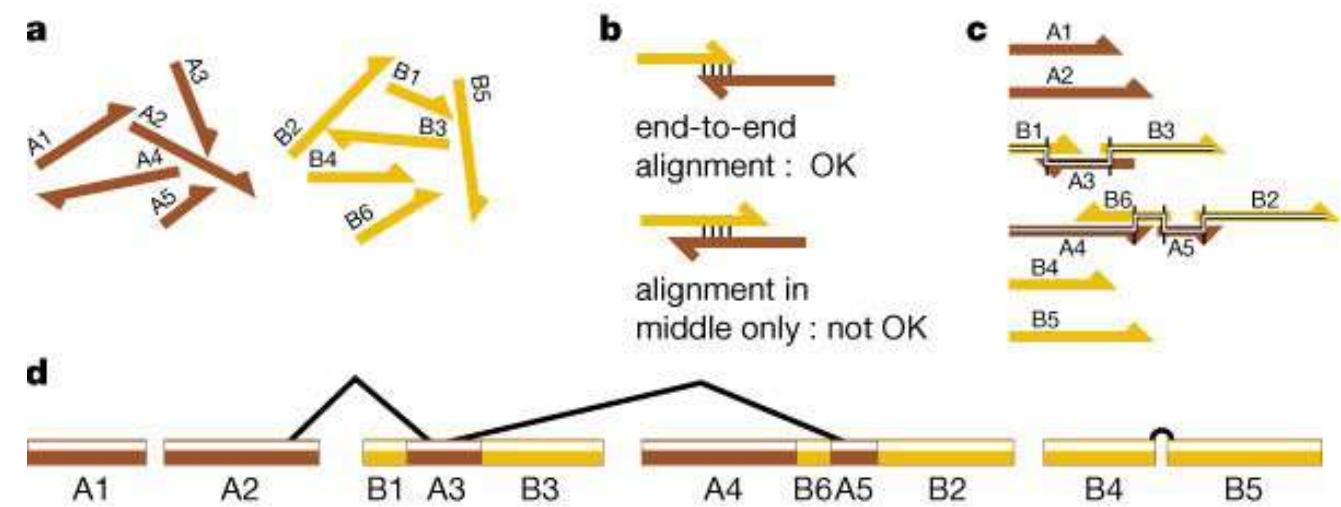
SÓ FALTA
MONTAR



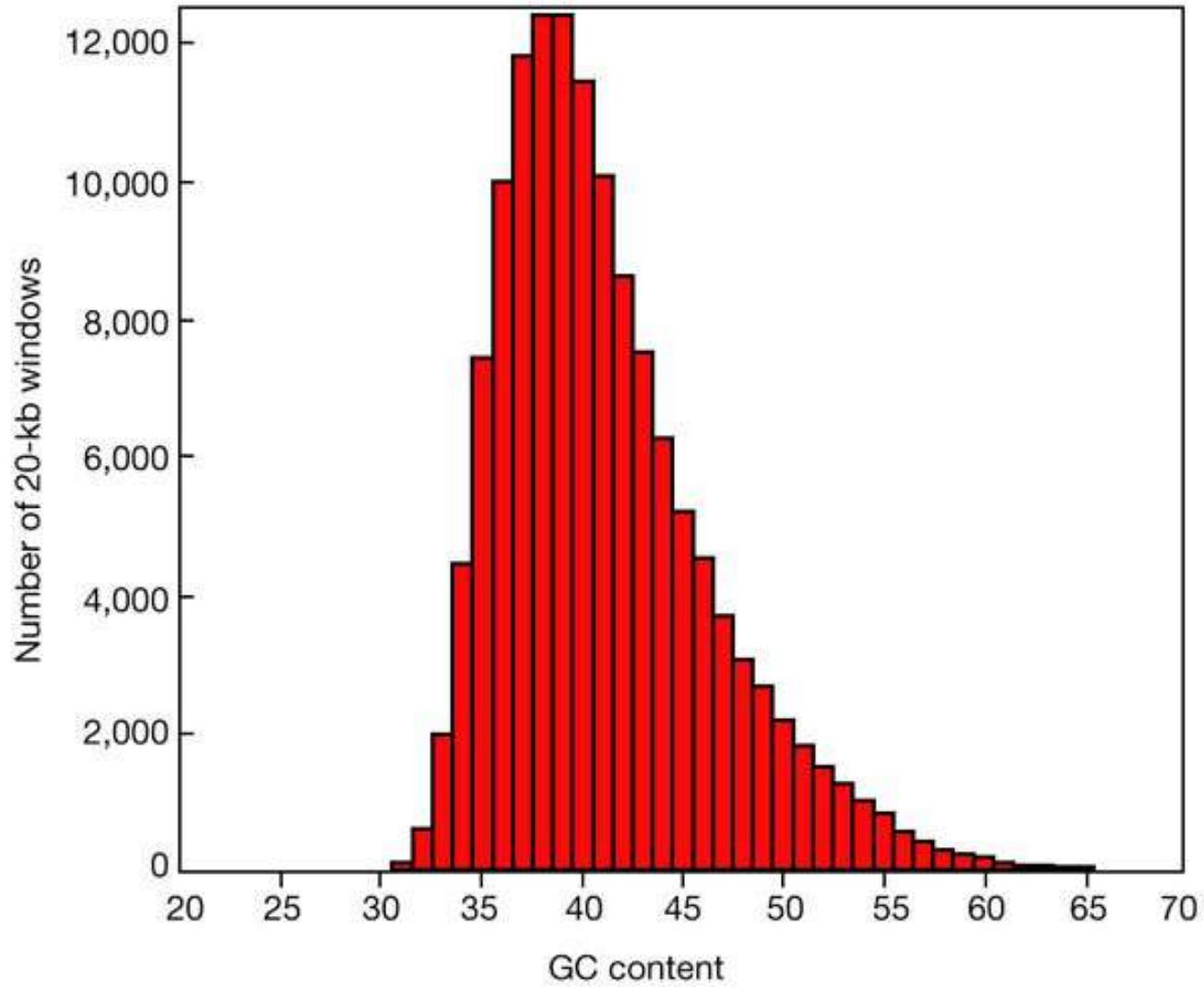
HIGH THROUGHPUT GENOME SEQUENCE



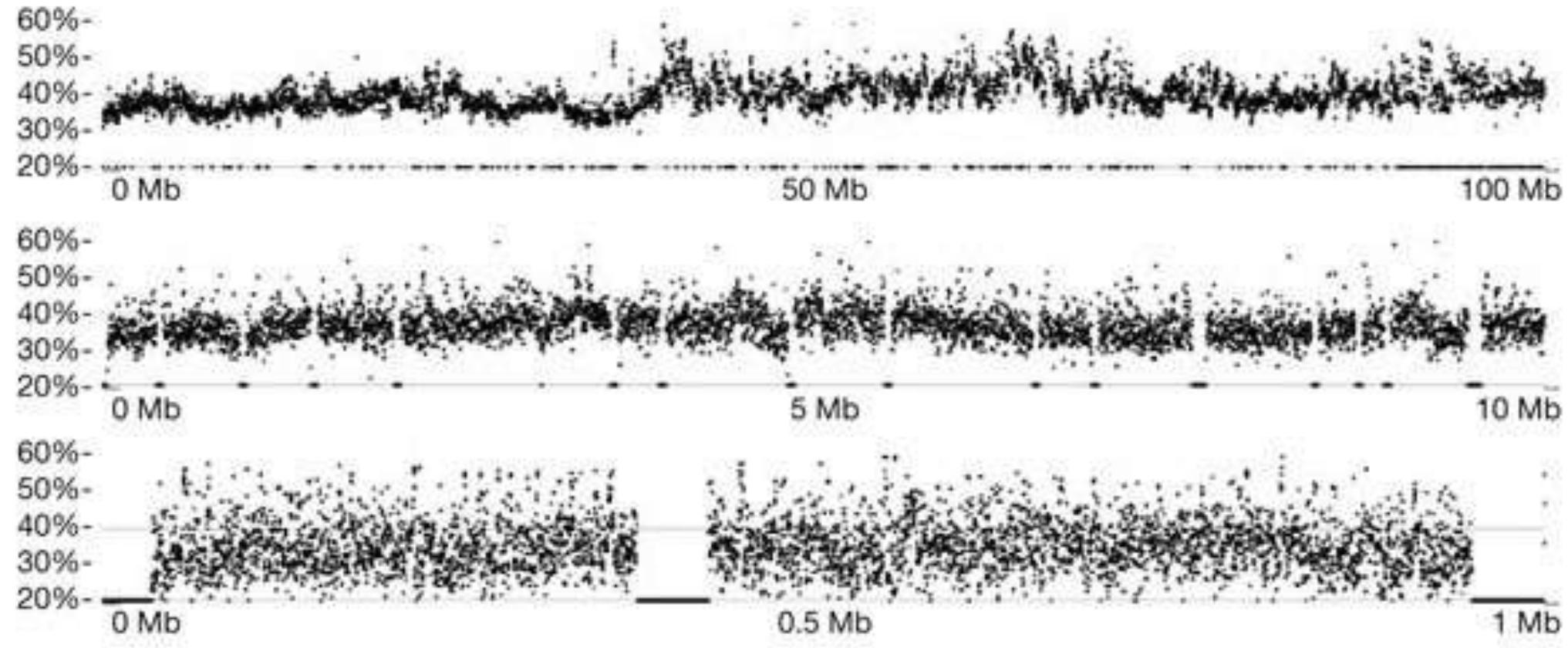
ASSEMBLING CLONES INTO THE DRAFT GENOME



CONTEÚDO GC

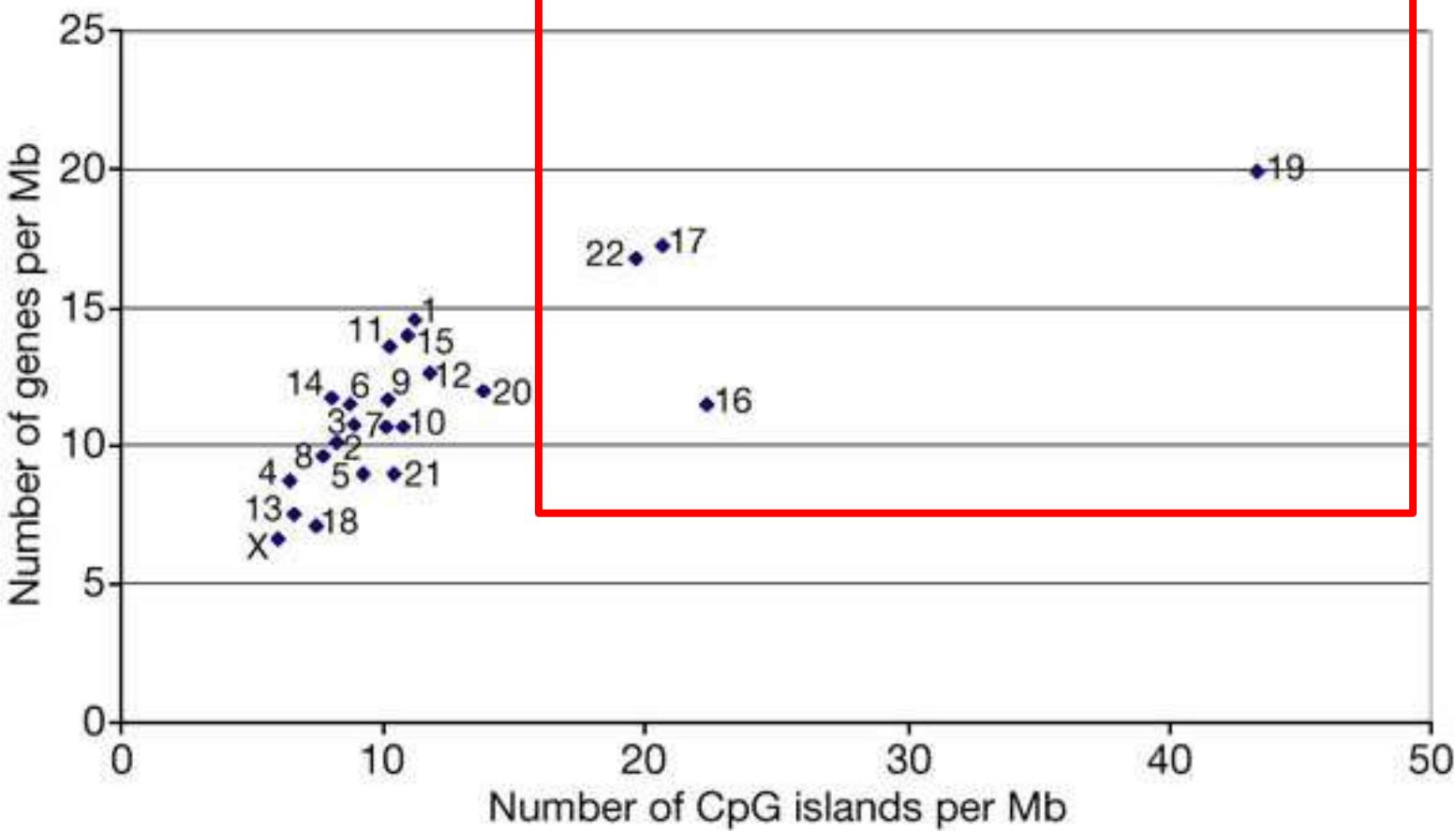


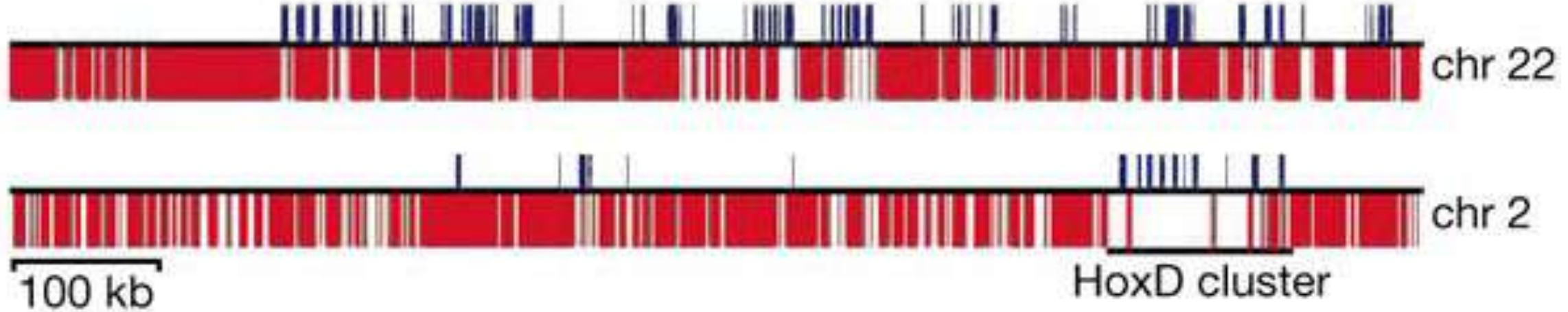
Variation in GC content at various scales.



100-Mb (20-kb windows).
10 Mb (2-kb windows).
1 Mb (200-bp windows)

NUMBER OF CPG ISLANDS PER MB FOR EACH CHROMOSOME

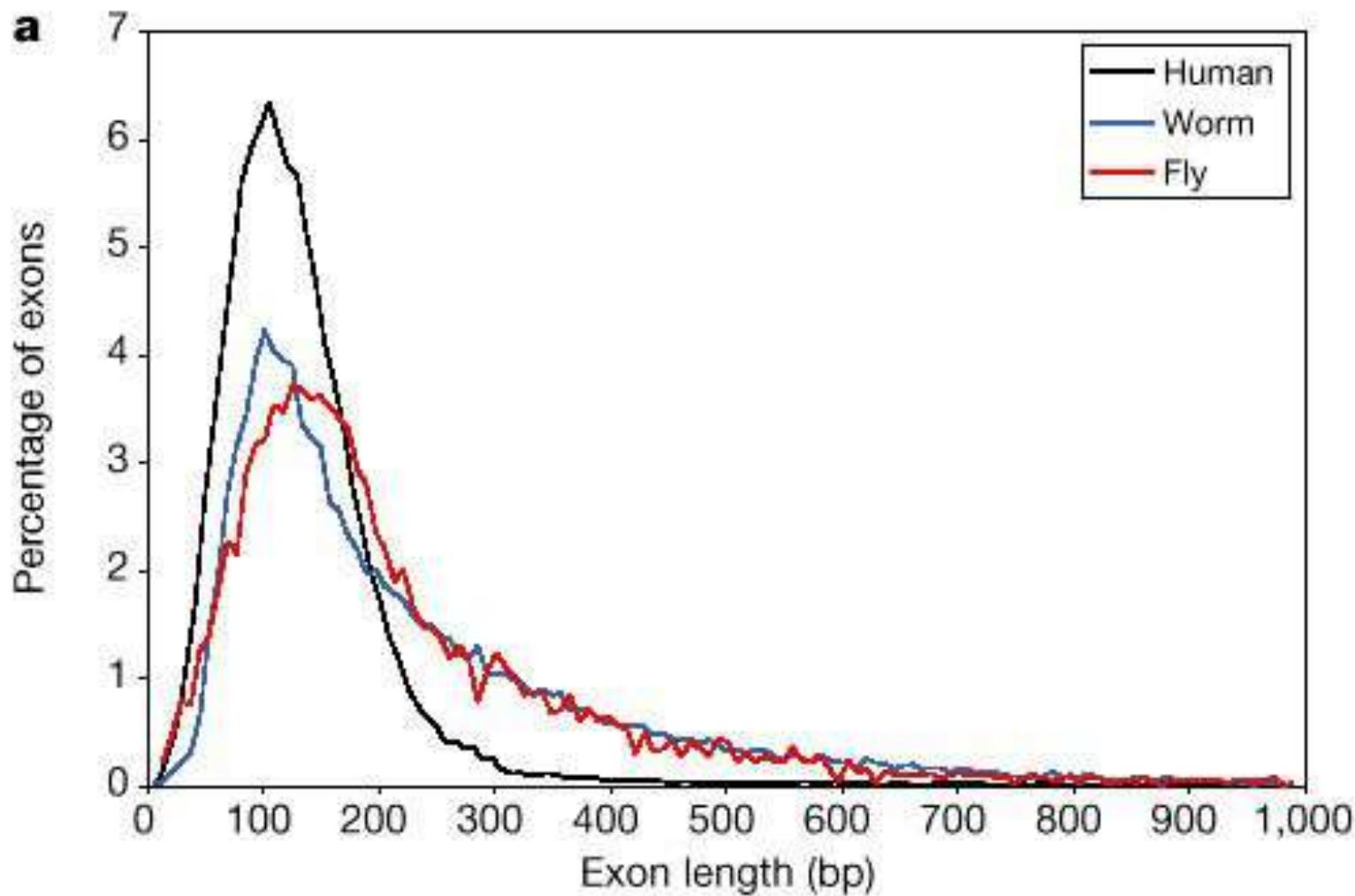




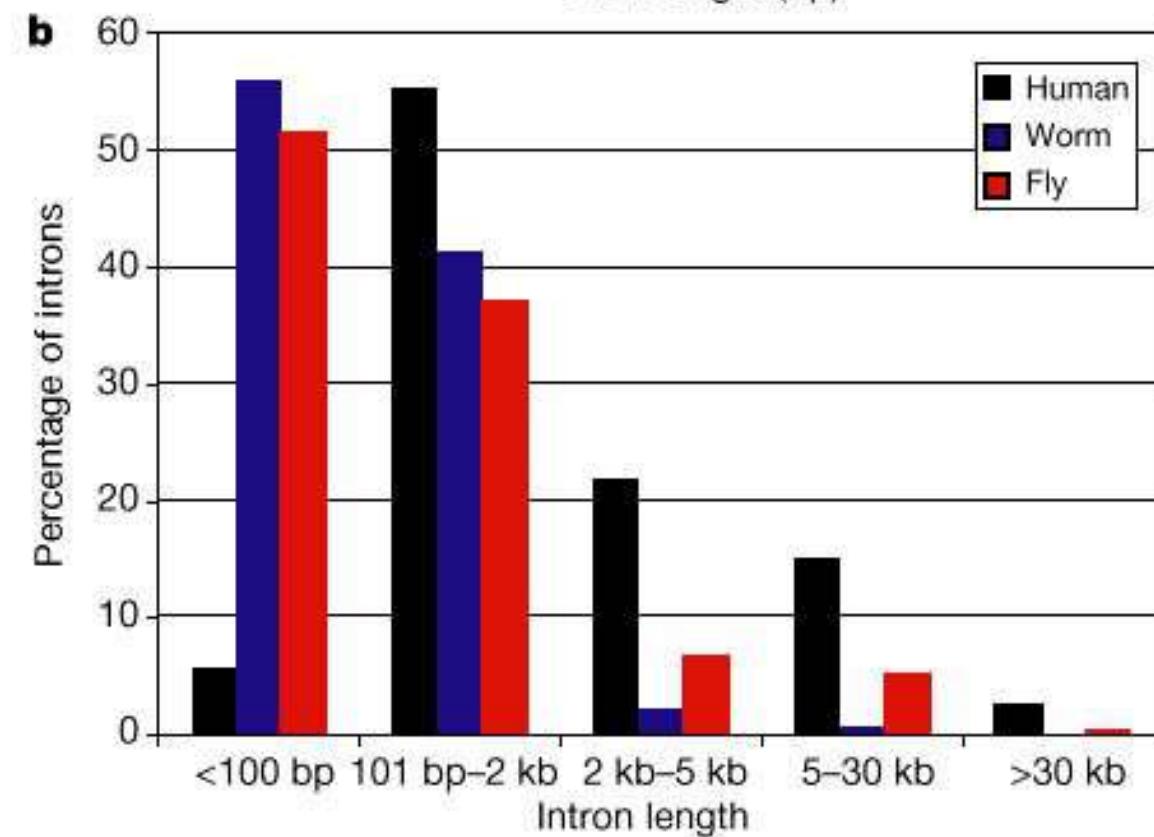
INTERSPERSED REPEATS

interspersed repeats → Vermelho
Exons → Azul

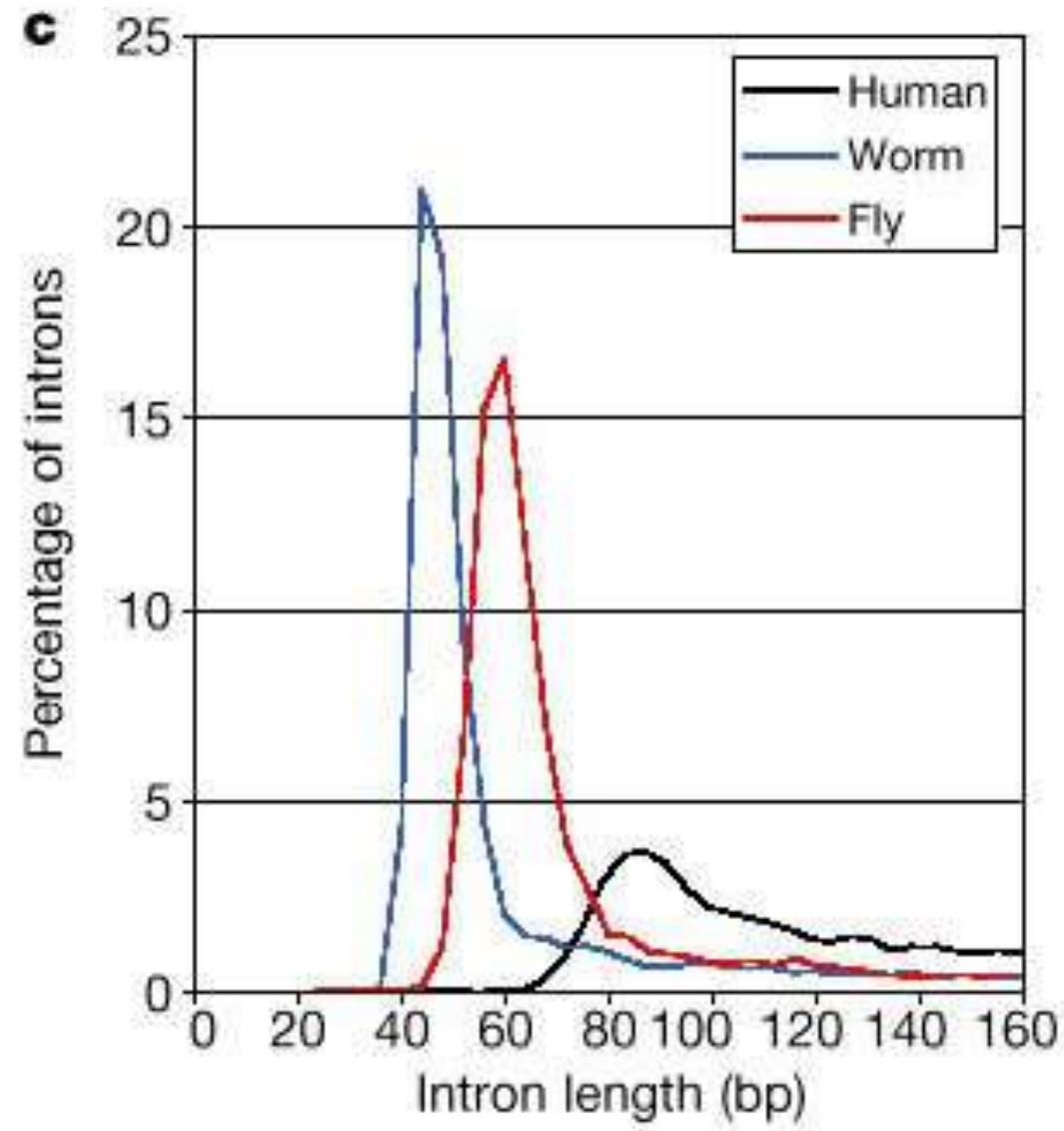
EXONS



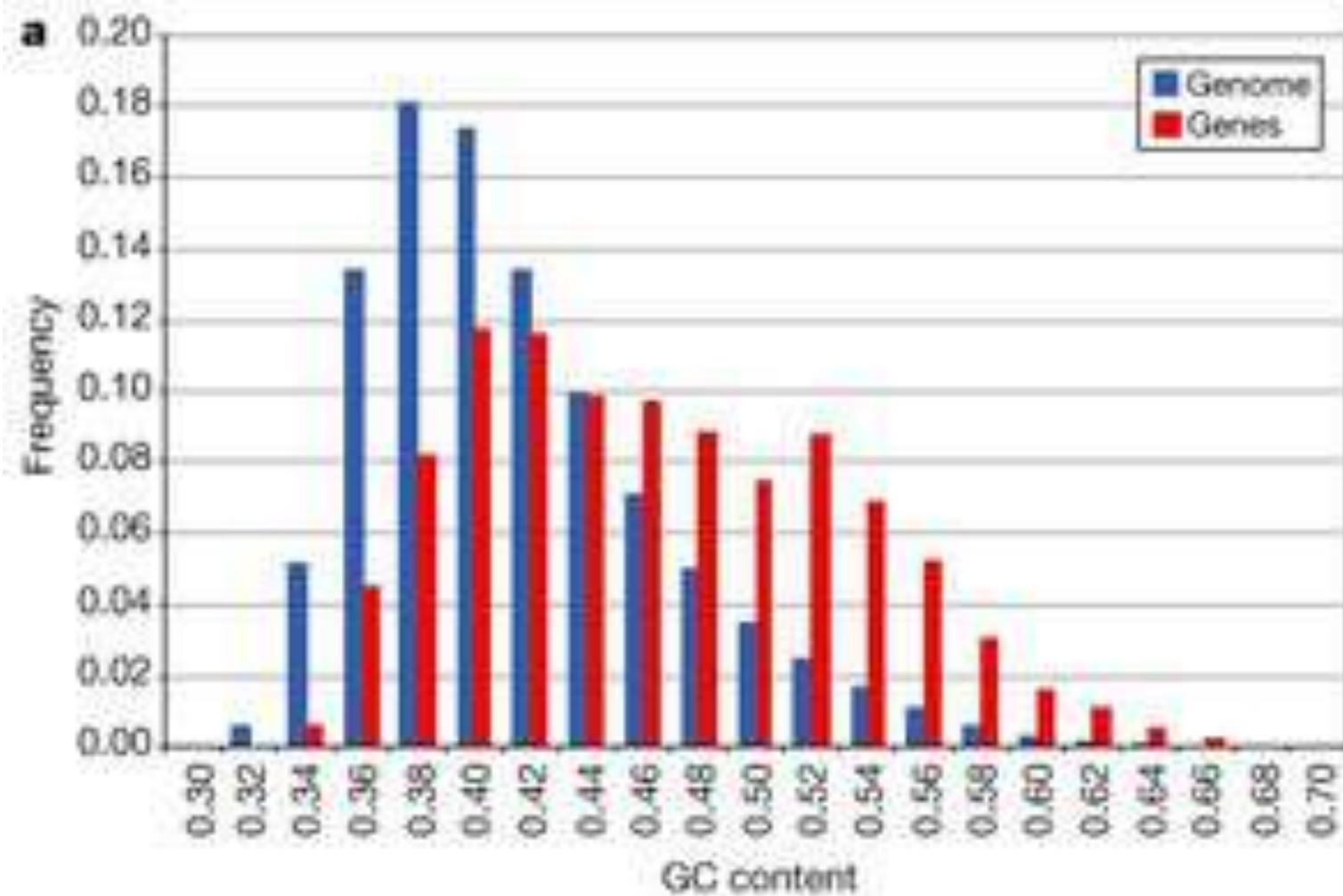
INTRONS



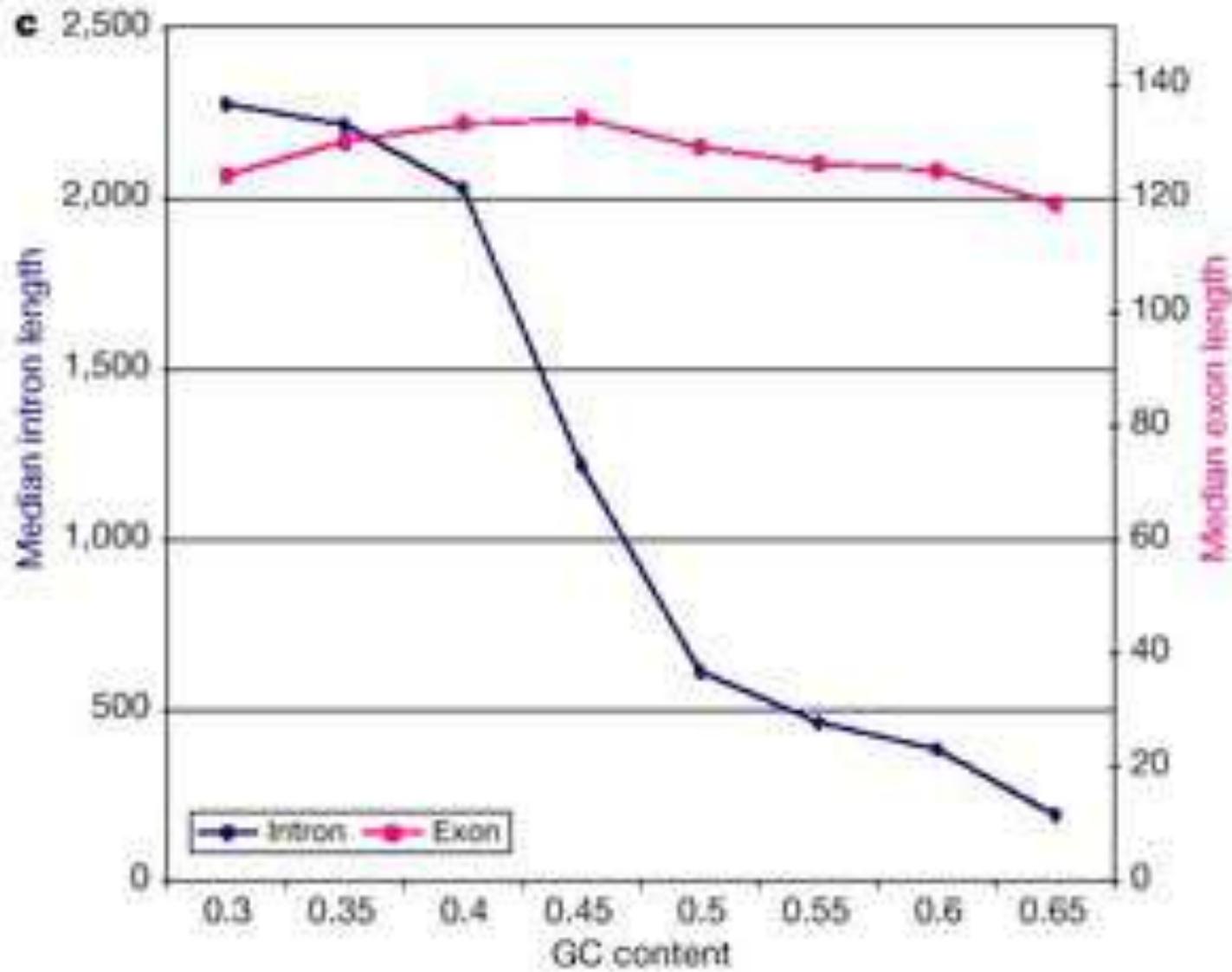
INTRONS



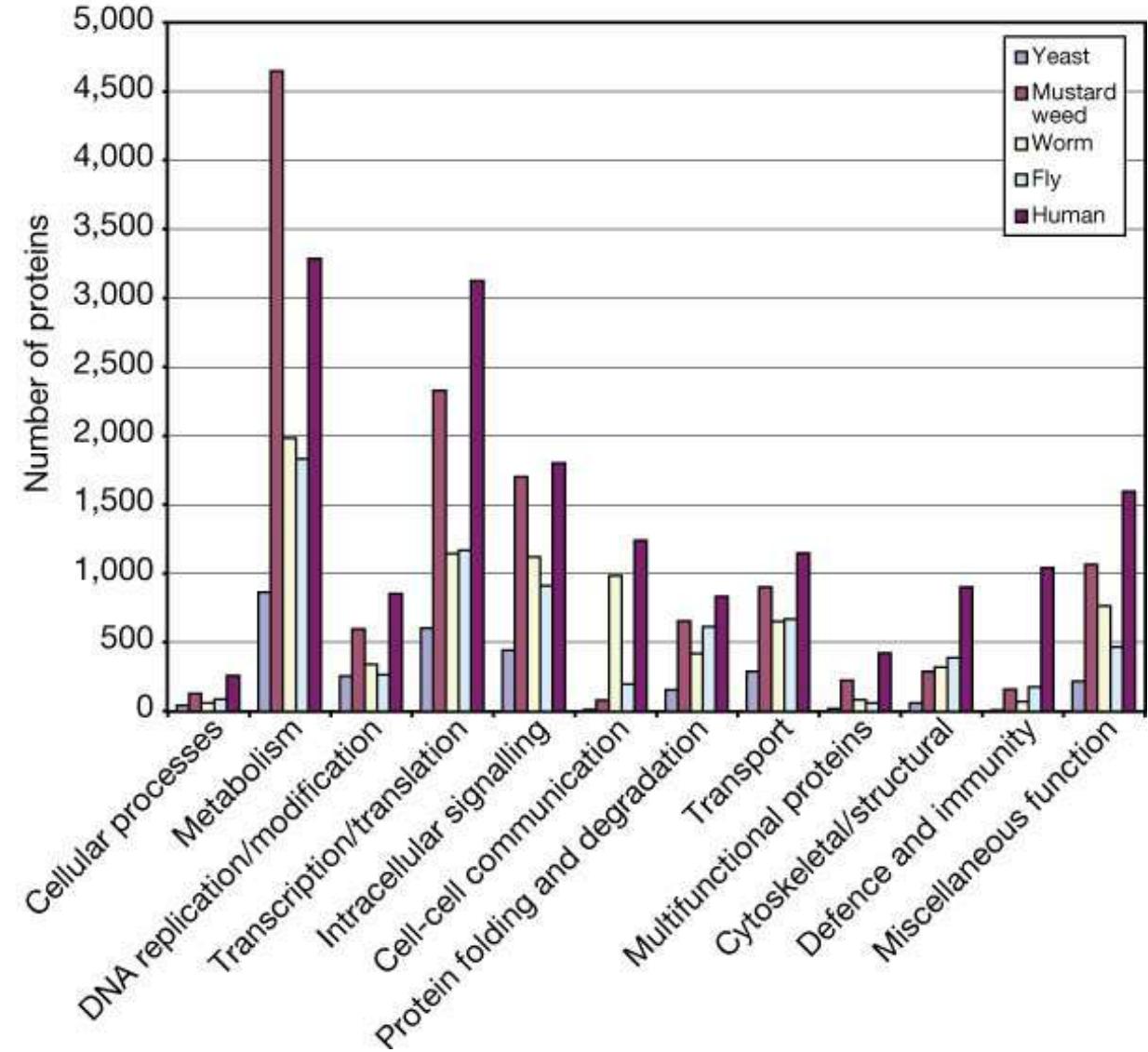
GC% GENOMA GENE



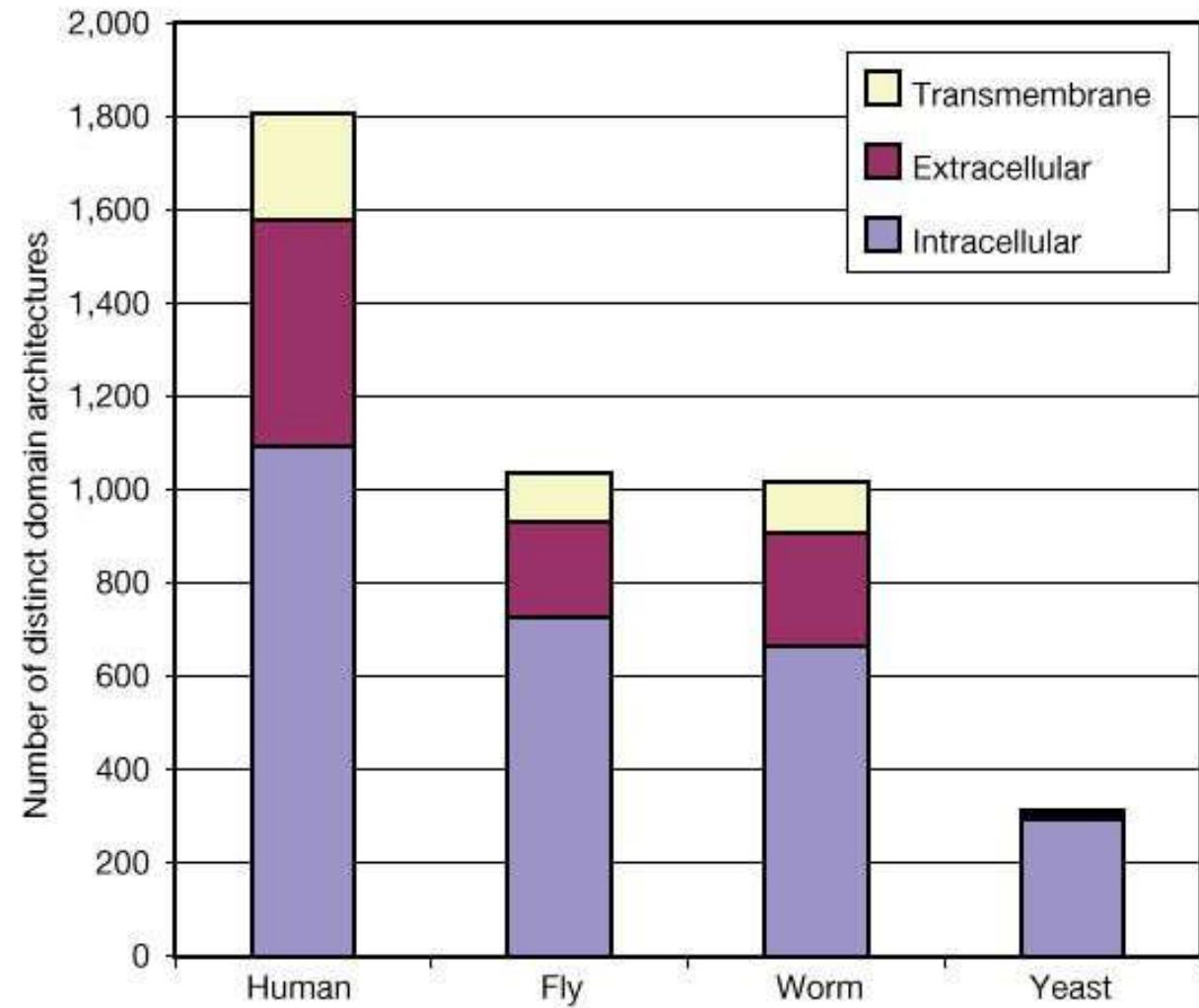
GC% EXON INTRON



FUNÇÕES DAS PROTEÍNAS

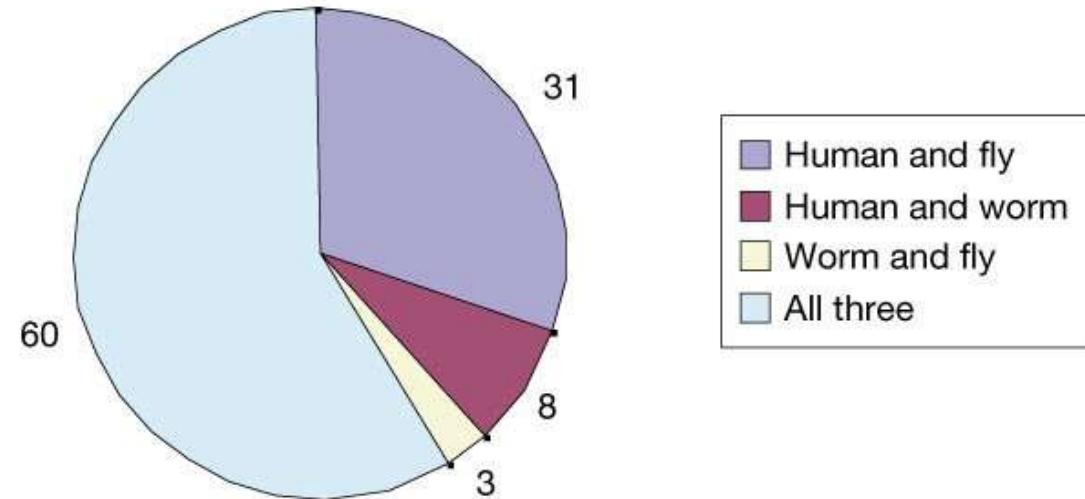


ARQUITETURA DE PROTEÍNAS

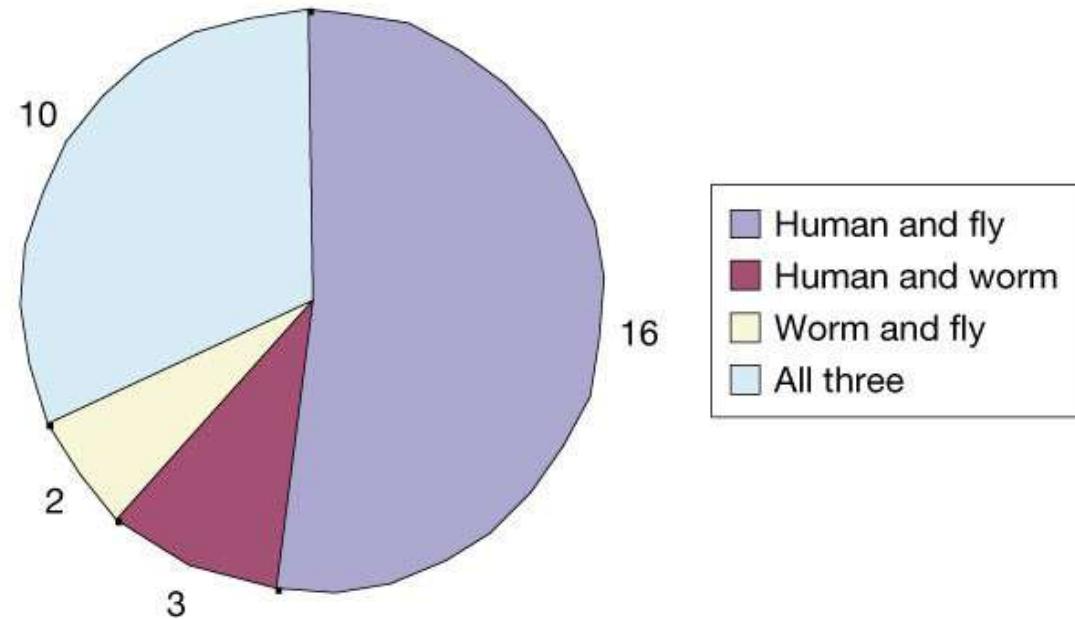


ARQUITETURA DE PROTEÍNAS

a Conserved domain architectures in chromatin proteins

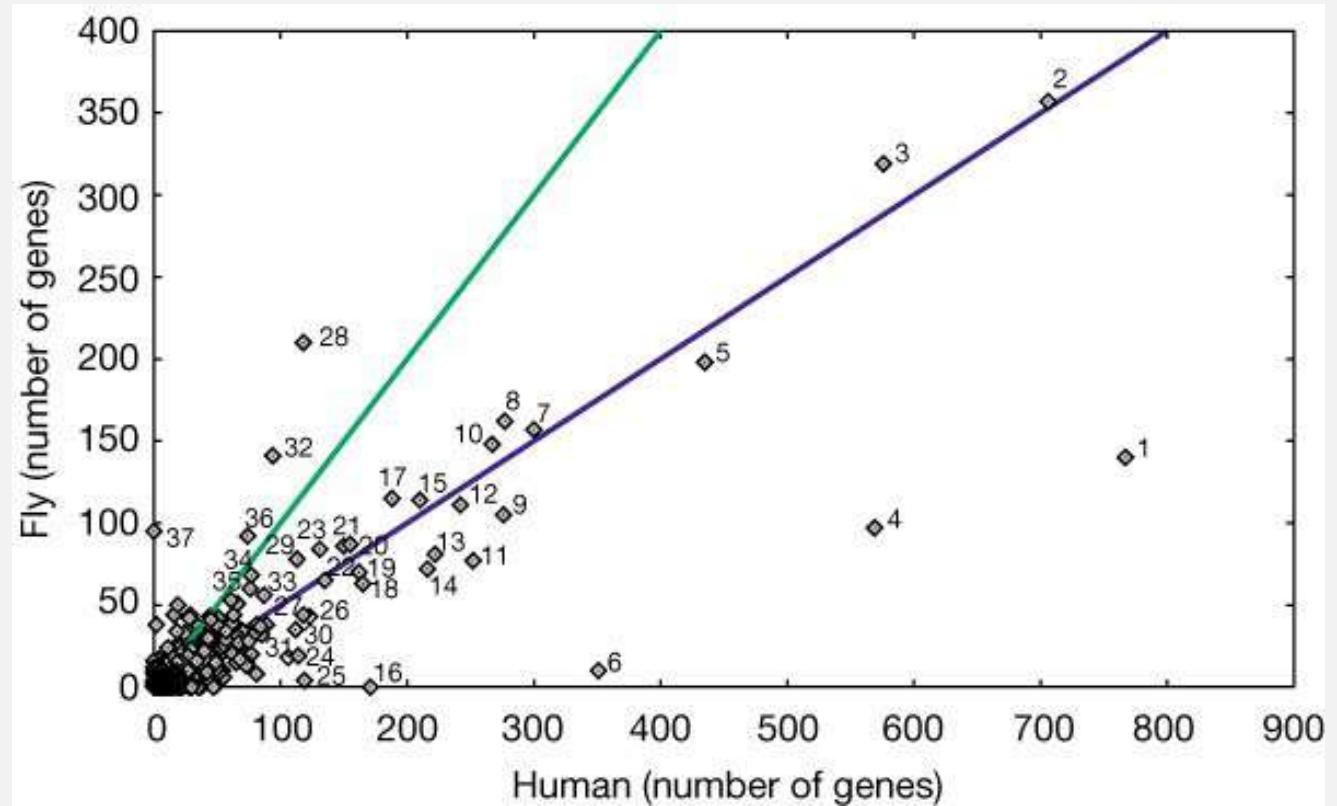


b Conserved domain architectures in apoptotic proteins

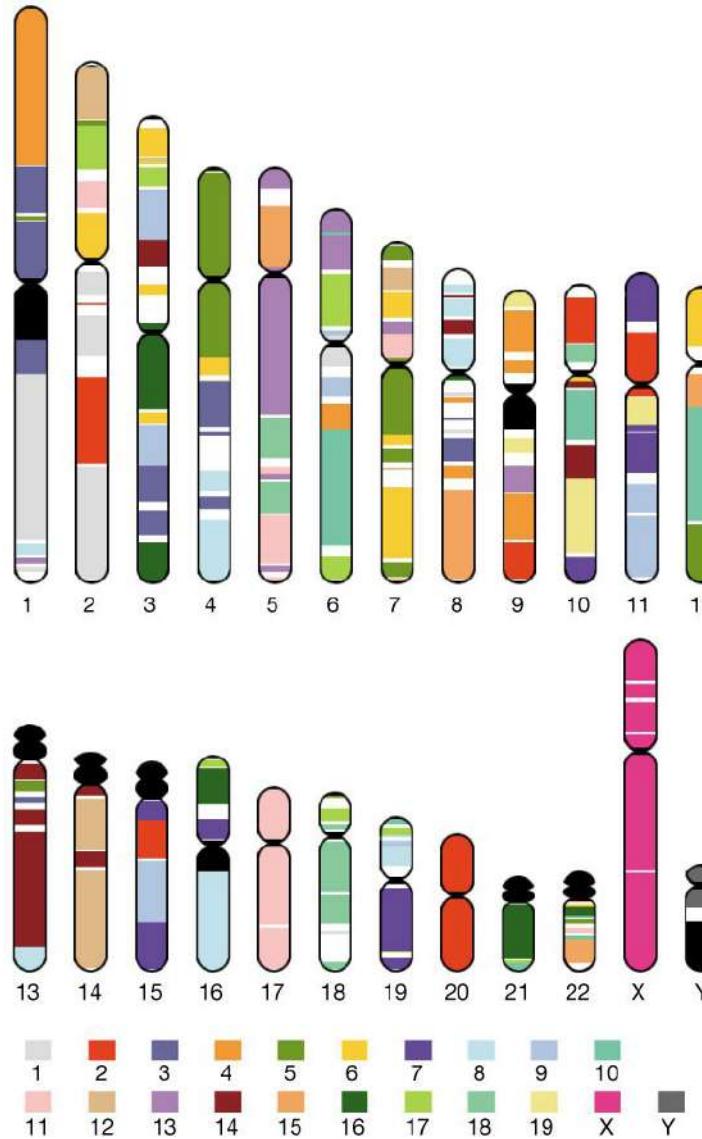


EXPANSÃO DAS FAMÍLIAS DE PROTEÍNAS

- (1) immunoglobulin domain
- (4) rhodopsin-like GPCR superfamily [IPR000276]
- (6) reverse transcriptase (RNA-dependent DNA polymerase)
- (28) serine proteases, trypsin family [IPR001254]



CONSERVAÇÃO DE GENOMAS HUMANO X MICE

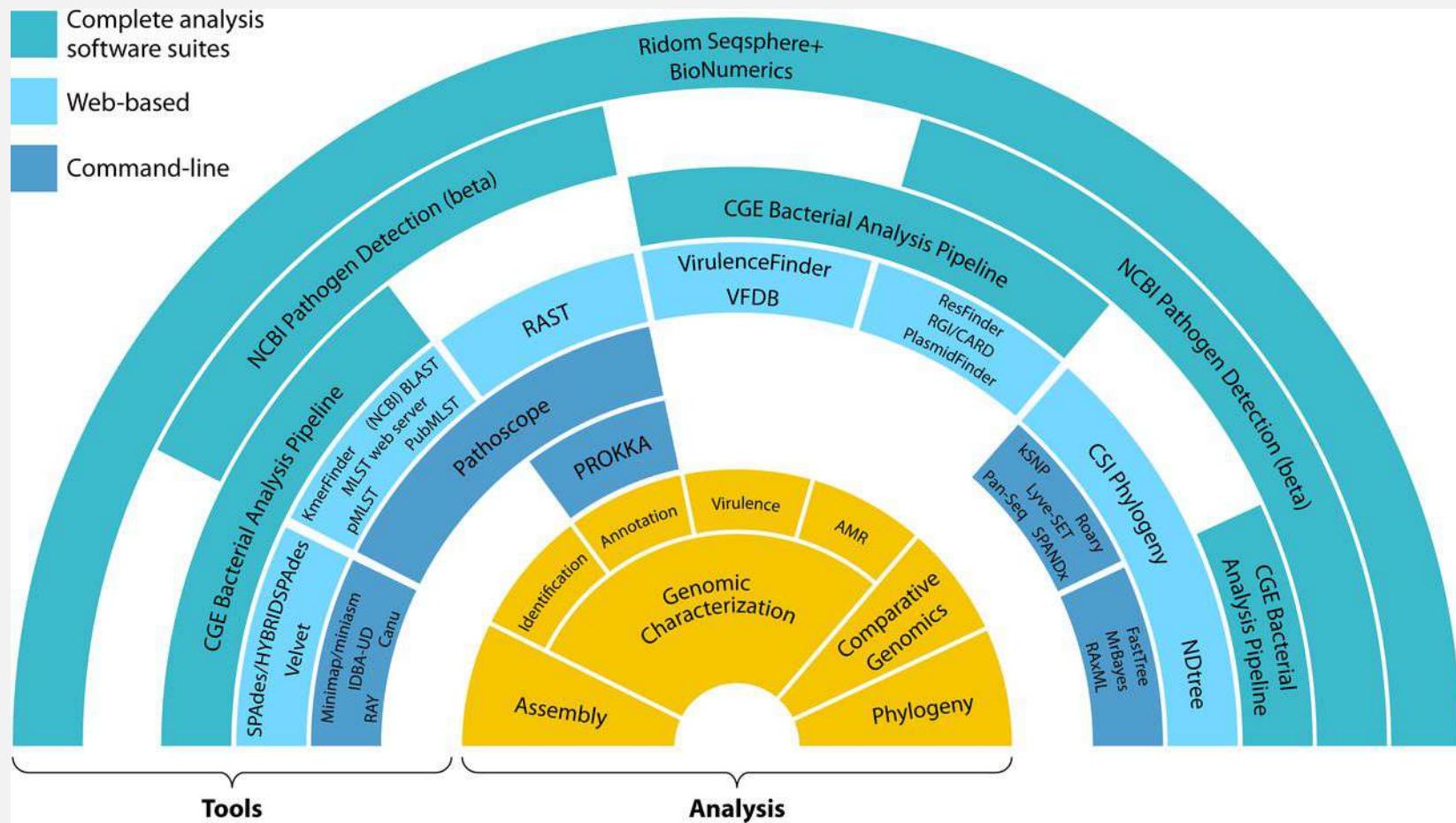


MONTAGEM DE GENOMA

Prof. Leandro Martins de Freitas, PhD
IMS/UFBA



WGS outbreak analysis tools.



Scott Quainoo et al. Clin. Microbiol. Rev. 2017;
doi:10.1128/CMR.00016-17



GENOMA

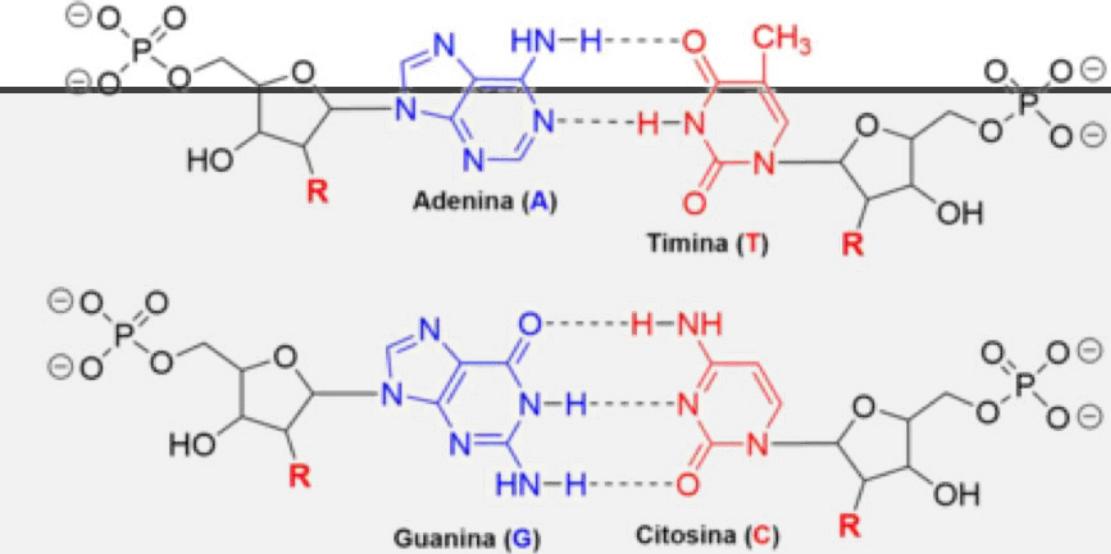
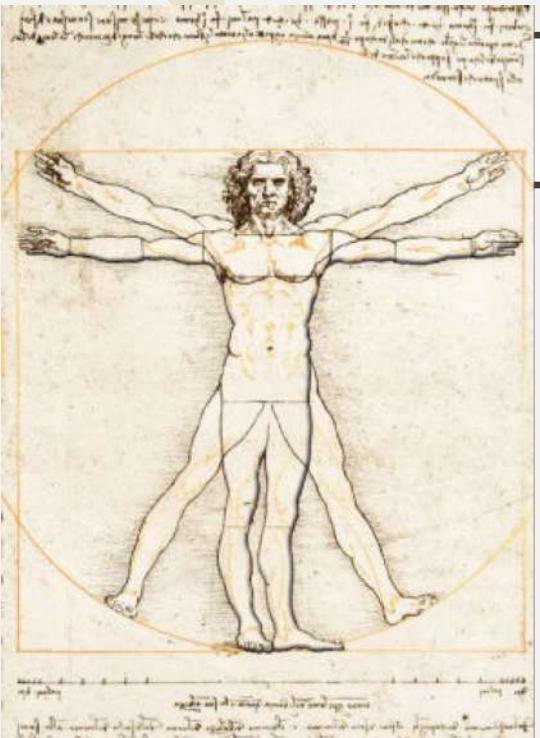
Genoma é todo o conteúdo de DNA presente em um organismo

Todos os cromossomos

Todos os genes

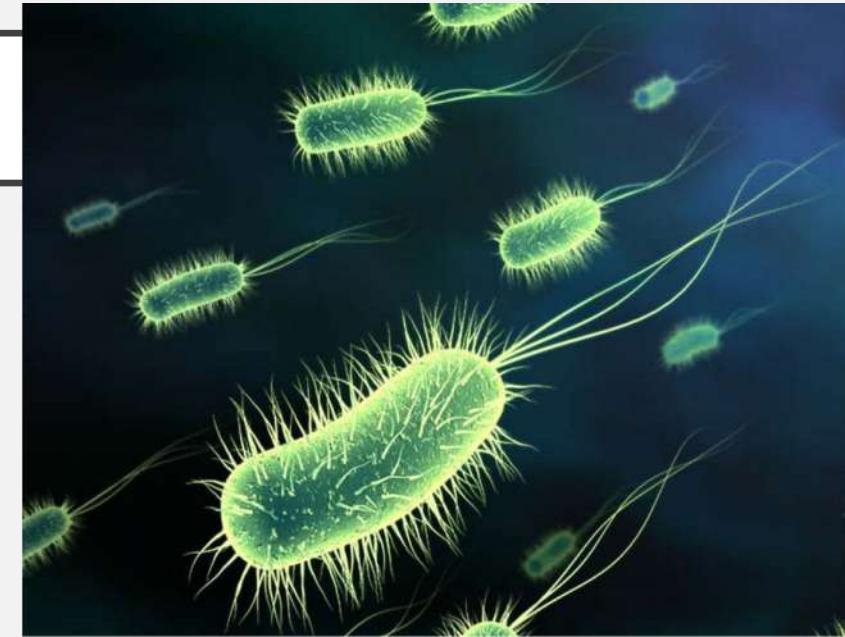
Elementos regulatórios e estruturais

Genômica: estudo em larga escala de genes, elementos regulatórios e estruturais contidos na sequência de DNA de um organismo



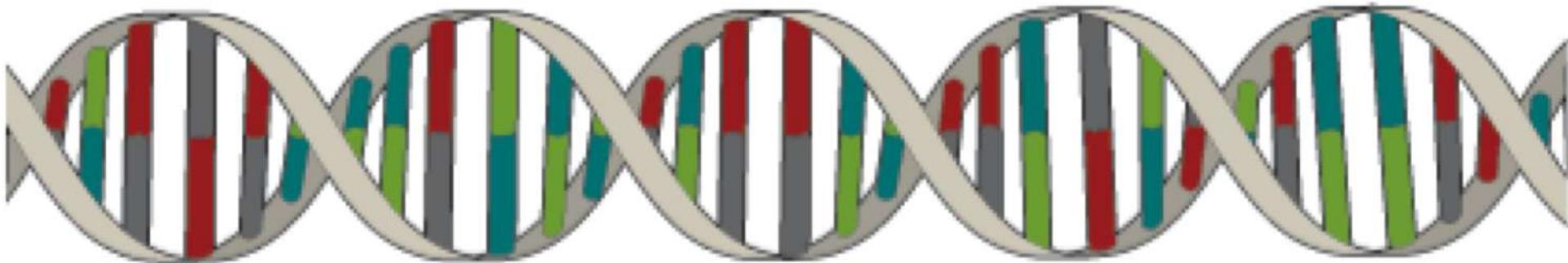
3 bilhões de bases

22 mil genes

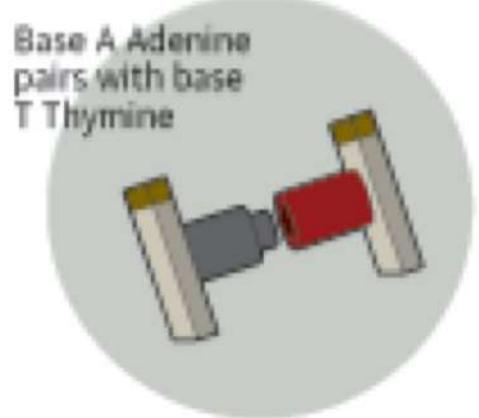
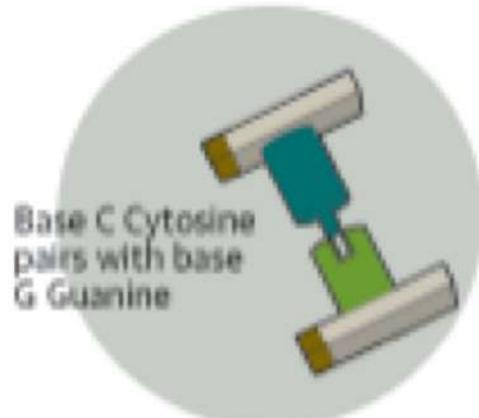
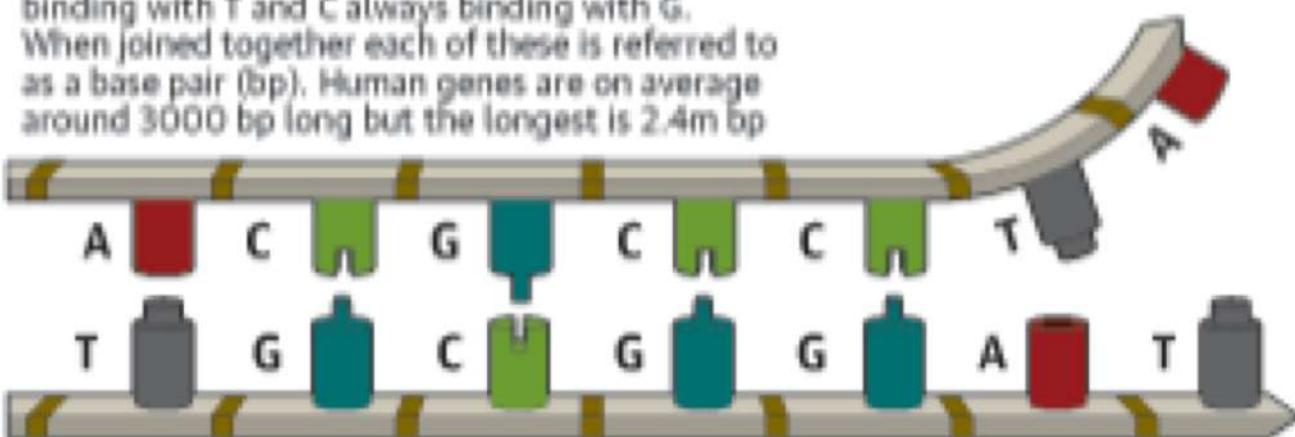


3 milhões de bases

3 mil genes



The two DNA strands that make up the double helix are joined by bonds between each of the bases. These complementary with A always binding with T and C always binding with G. When joined together each of these is referred to as a base pair (bp). Human genes are on average around 3000 bp long but the longest is 2.4m bp



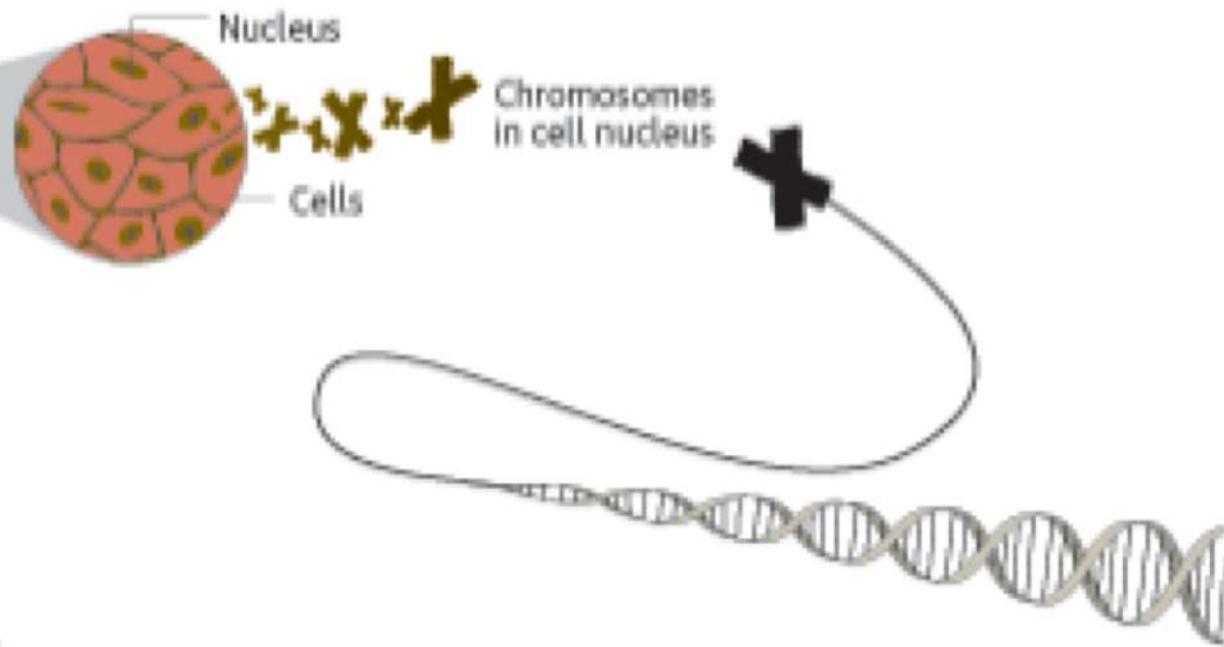
Different species vary greatly in the size of their genomes and the number of genes they have. There is a trend for increasing gene number with increasing

CHARTING THE HUMAN GENOME MAP

All humans share

99.8%

of their genetic material.



98%

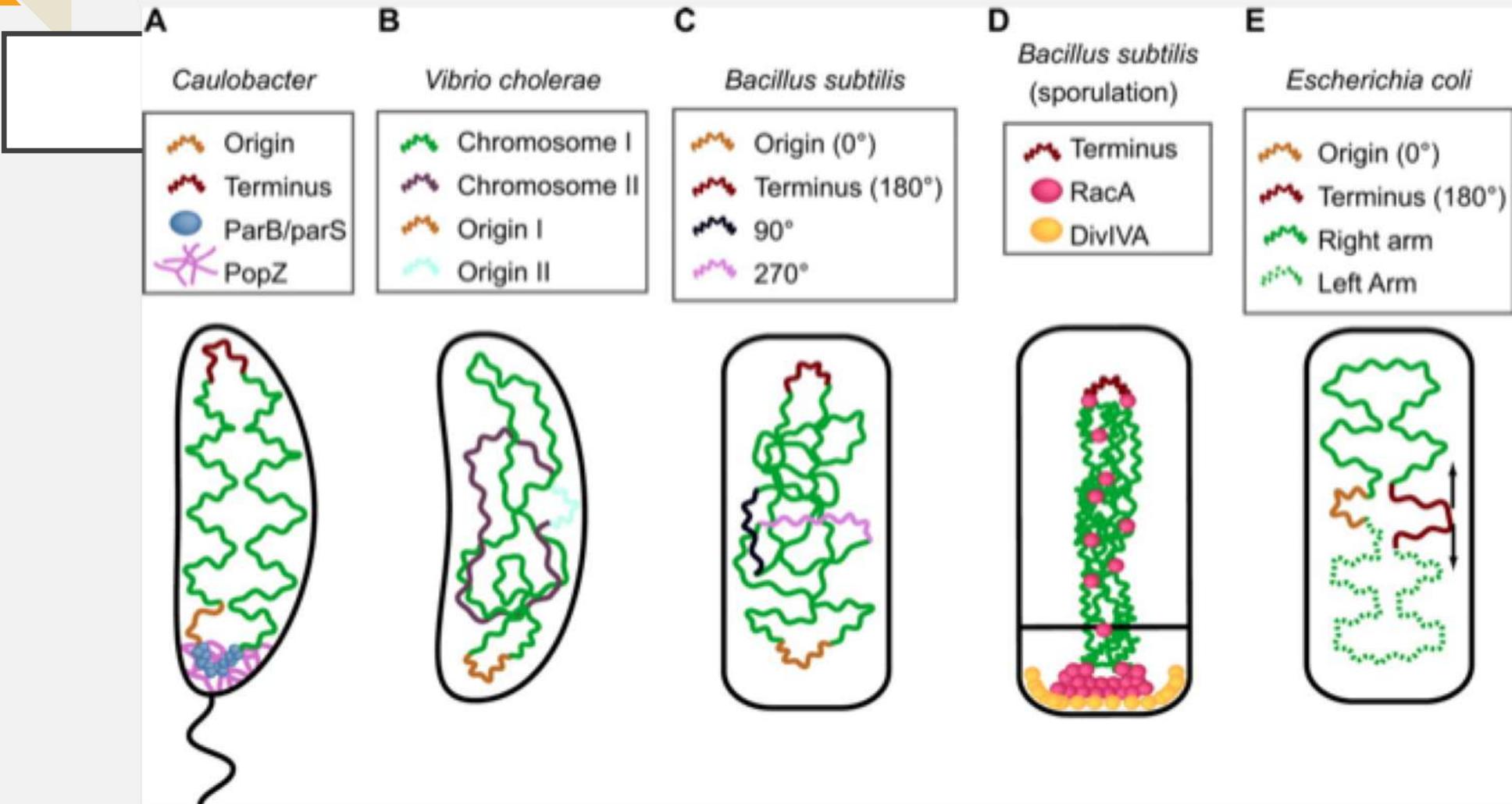
of human chromosomes are
"junk DNA" between genes. Large sections
are remnants of viral and bacterial infections

Different species vary greatly in the size of their genomes and the number of genes they have. There is a trend for increasing gene number with increasing body complexity, but there are numerous exceptions. Some species such as the mountain grasshopper, with 16.5 bn base pairs, have much larger genomes than us.

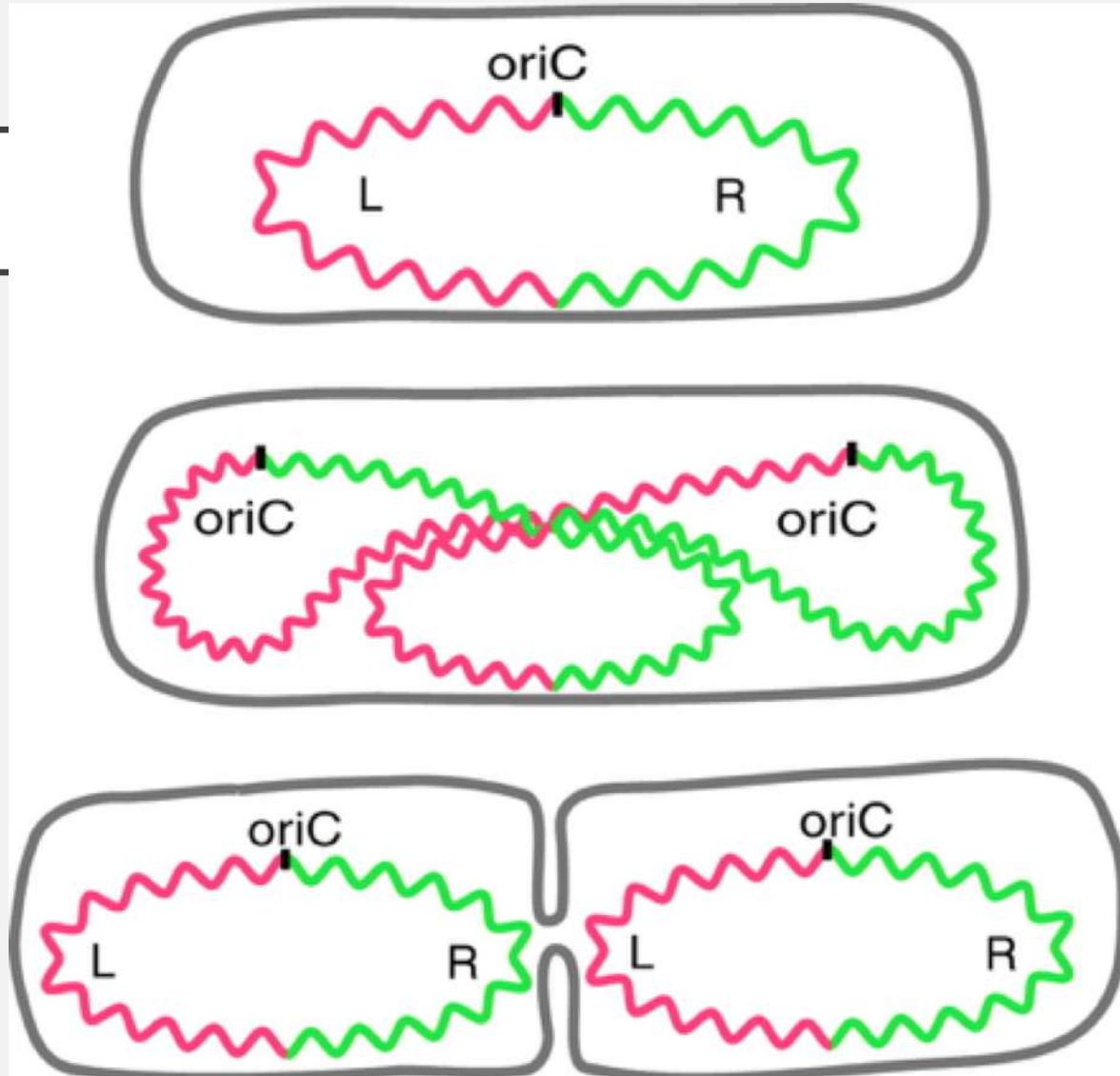


	Human	Thale cress	Earthworm	Fruit fly	Yeast
Genes	20-25,000	26,000	18,000	13,600	6,000
Base pairs	c.3 billion	117m	91m	116m	12m
Complexity	12 genes per million base pairs	221	197	117	500

Chromosome organization in model bacteria.

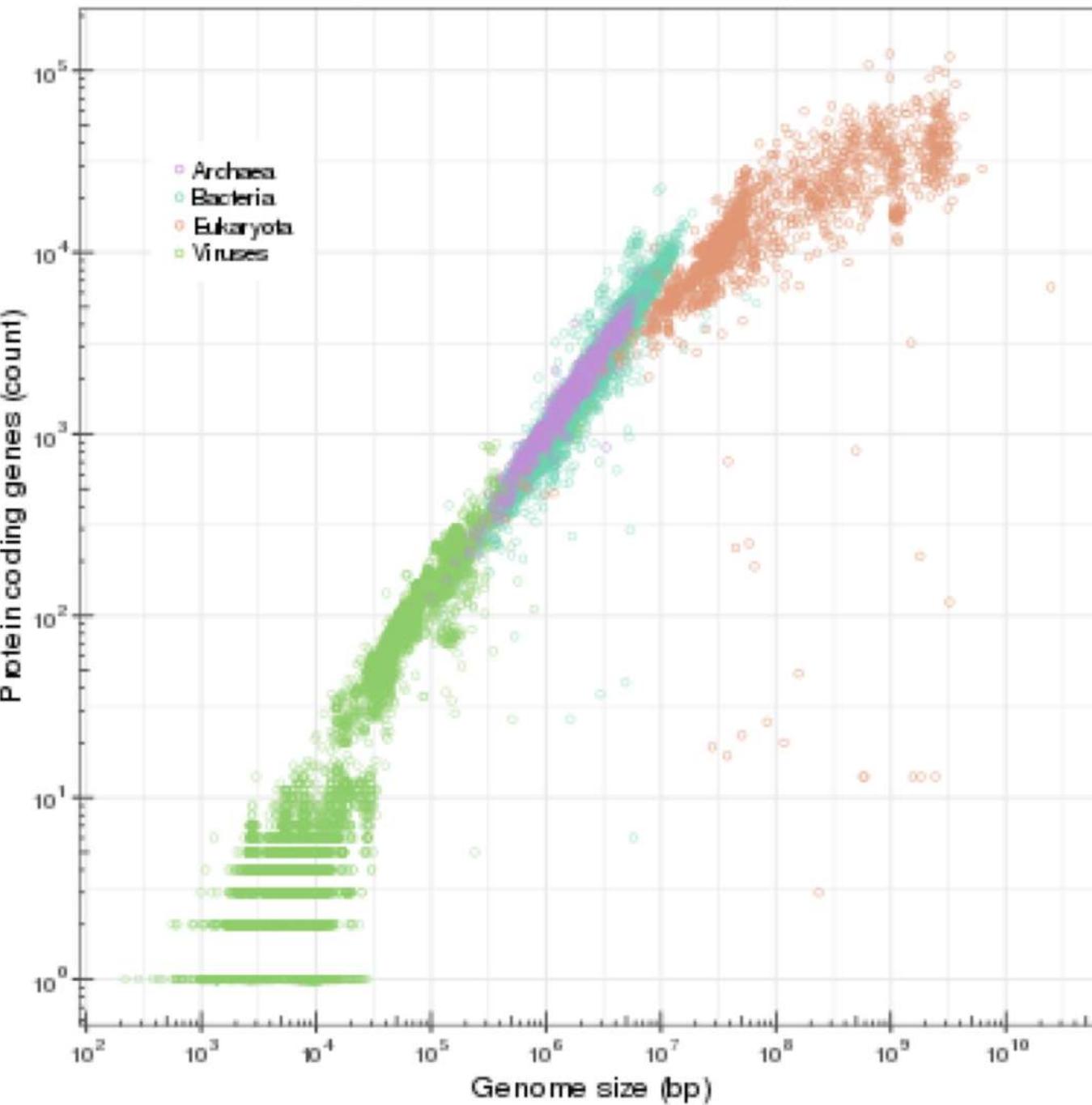


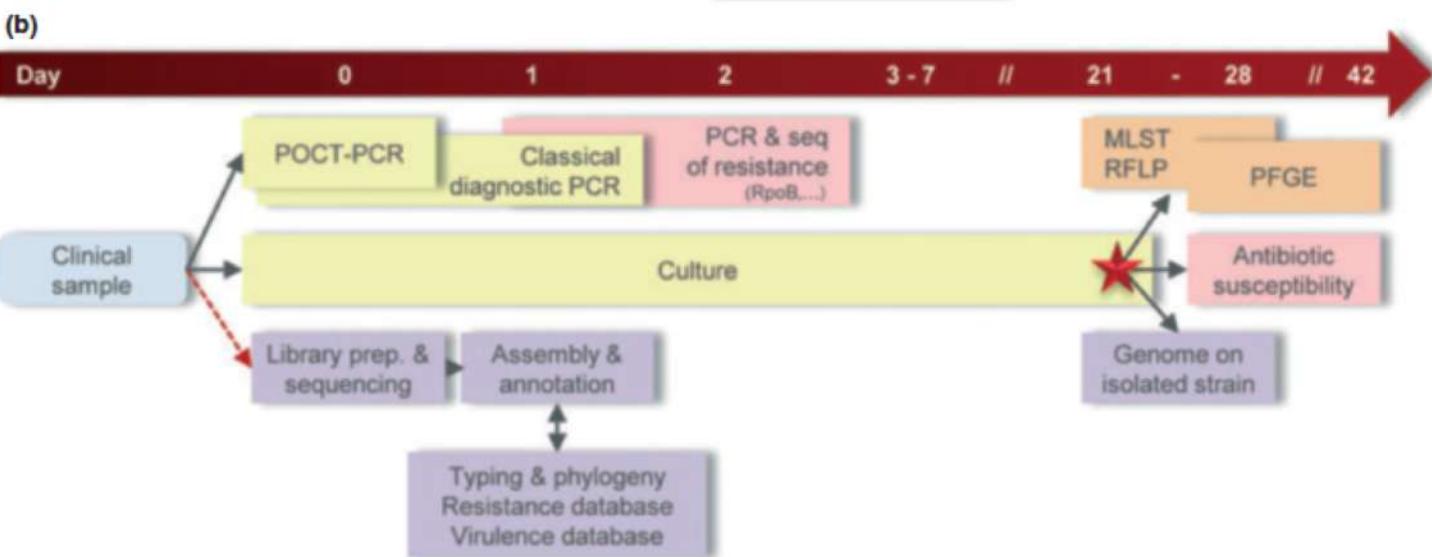
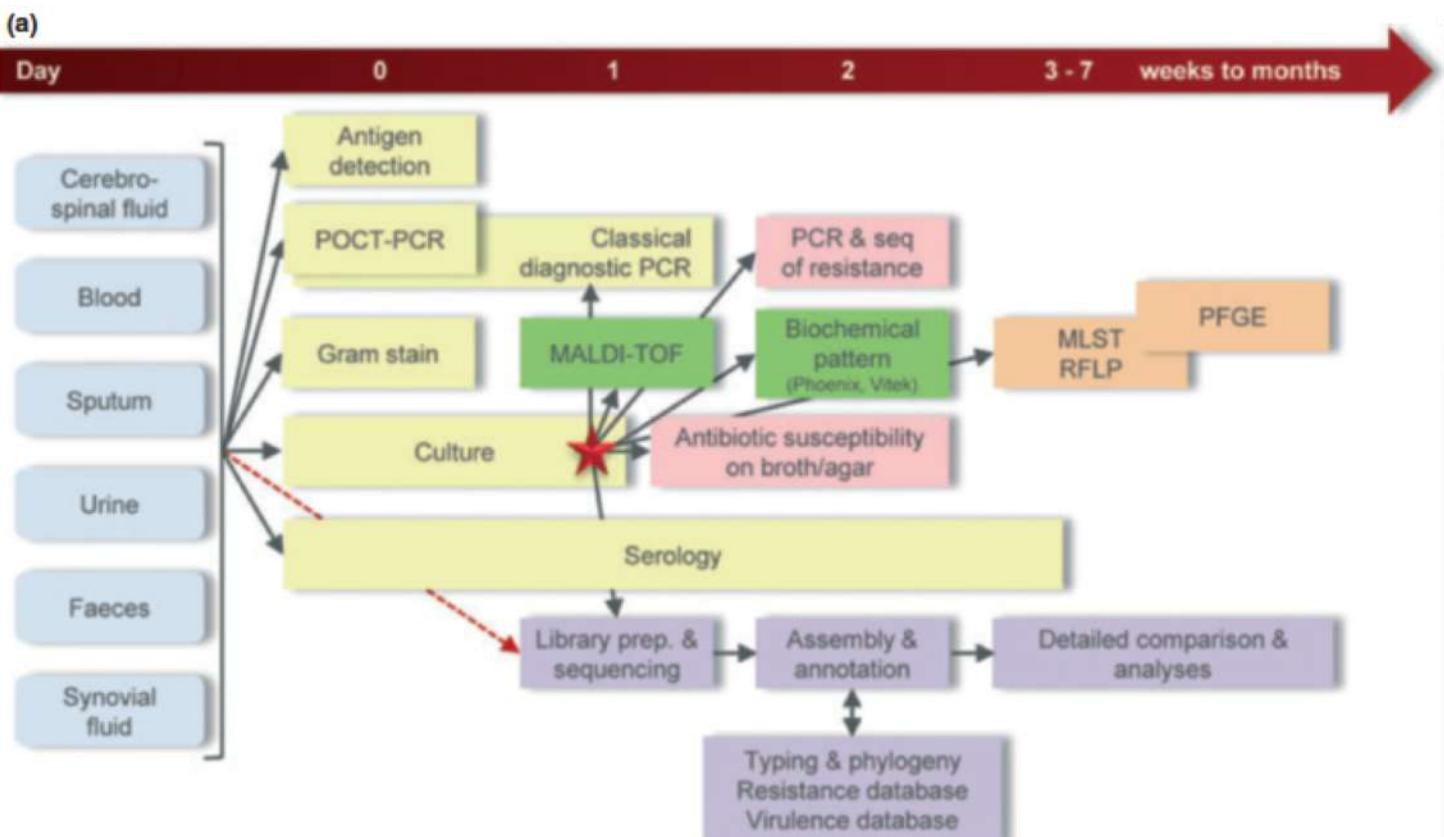
Chromosome organization in *E. coli*



Tamanho do genoma vs número de proteínas

Genome size vs. protein count across NCBI genomes





15 February 2001

nature

the
human
genome

www.nature.com

Nuclear fission

Five-dimensional
energy landscapes

Seafloor spreading

The view from under
the Arctic ice

Career prospects

Sequence creates new
opportunities

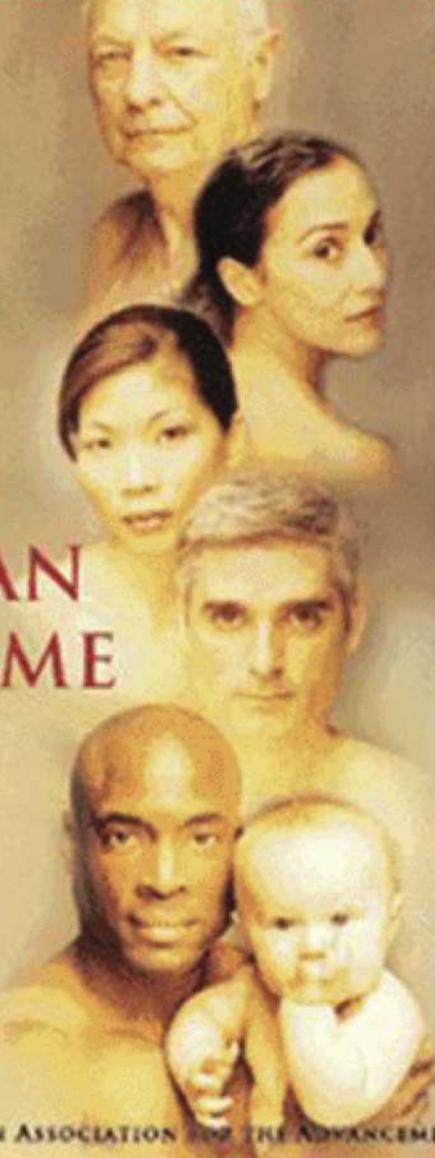
naturejobs
genomics special

Science

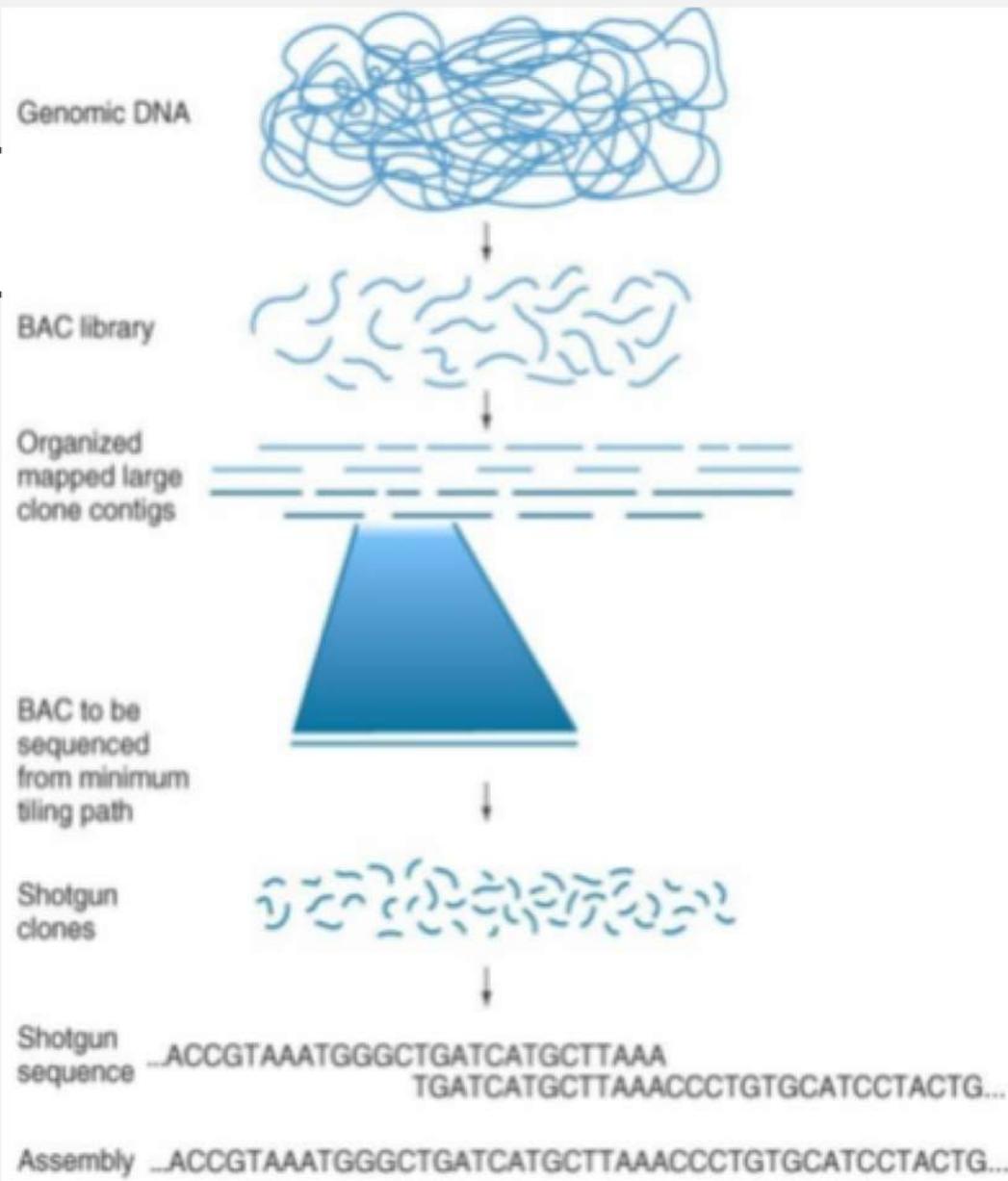
16 February 2001

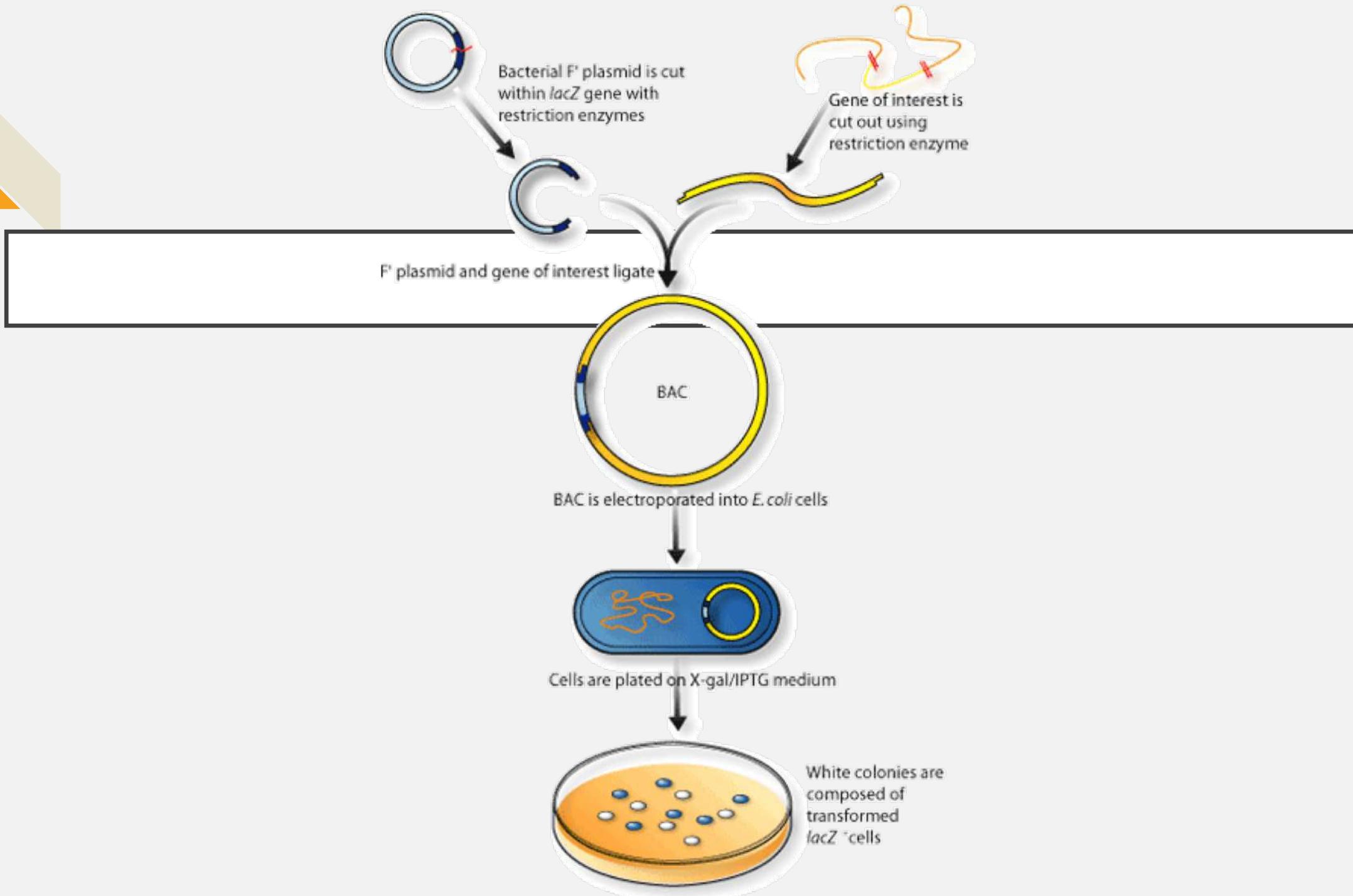
Vol. 291, No. 5507
Pages 1145-1434 \$9

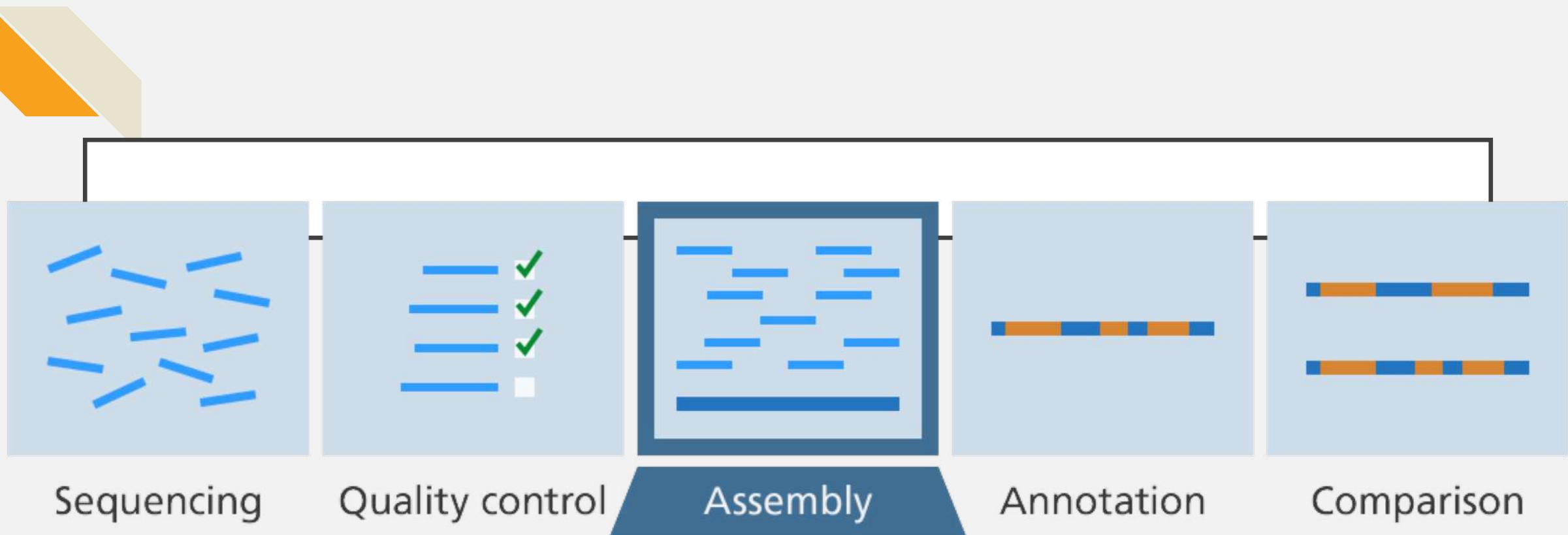
**THE
HUMAN
GENOME**



AMERICAN ASSOCIATION FOR THE ADVANCEMENT OF SCIENCE







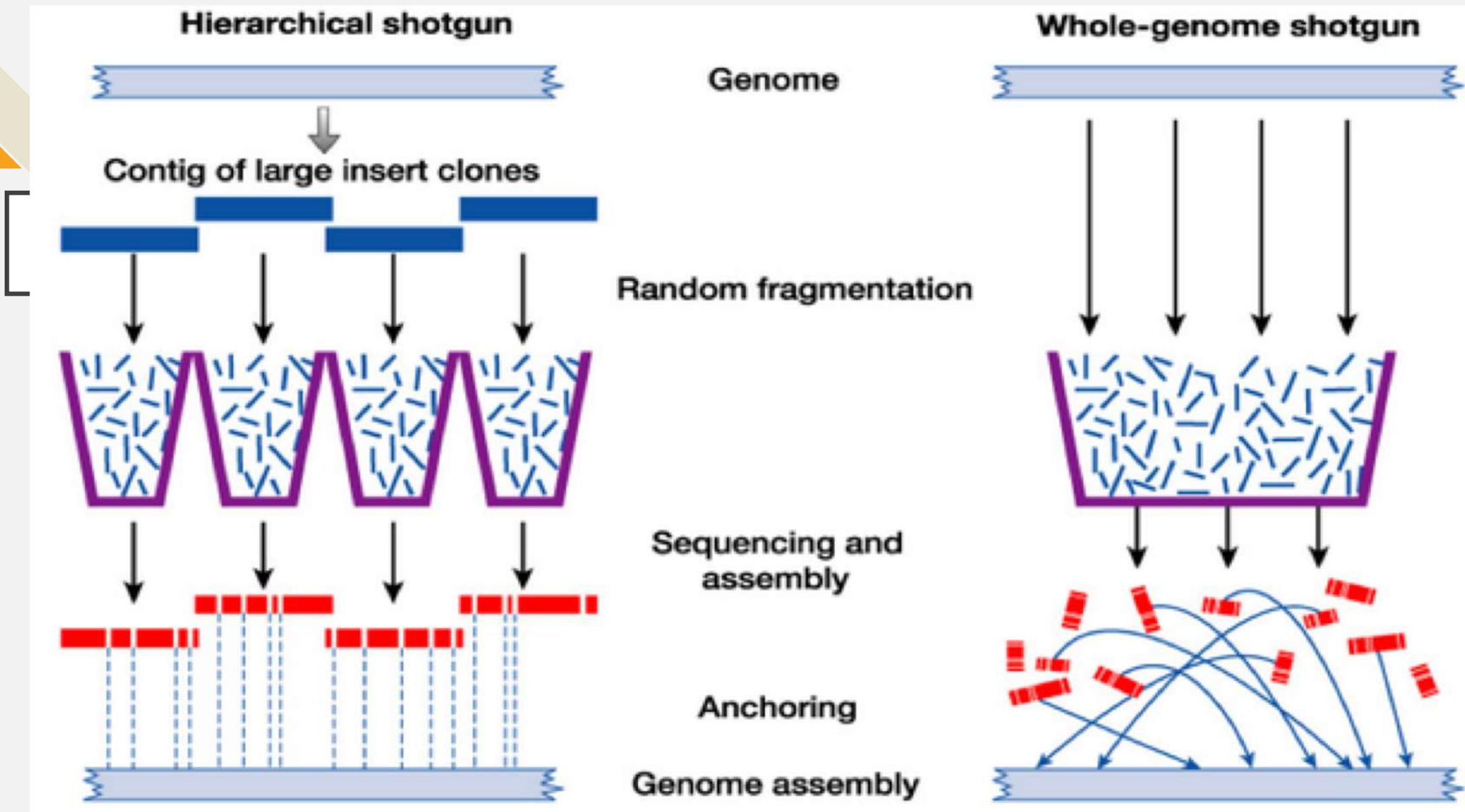


Figure 8-3 Human Molecular Genetics, 3/e. (© Garland Science 2004)



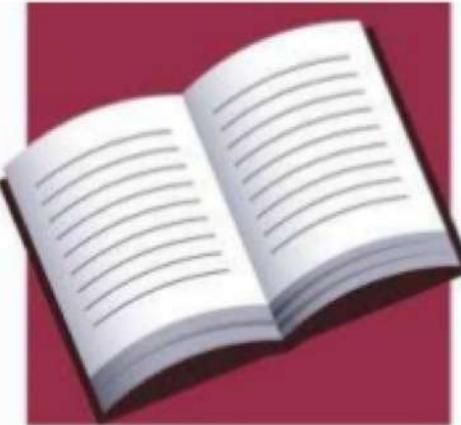
A single recipe is like...



A Gene:

One set of instructions for how to make one protein.

A recipe book is like...



A Chromosome:

Thousands of sets of instructions for how to make thousands of proteins

Two copies of 23 recipe books is like...

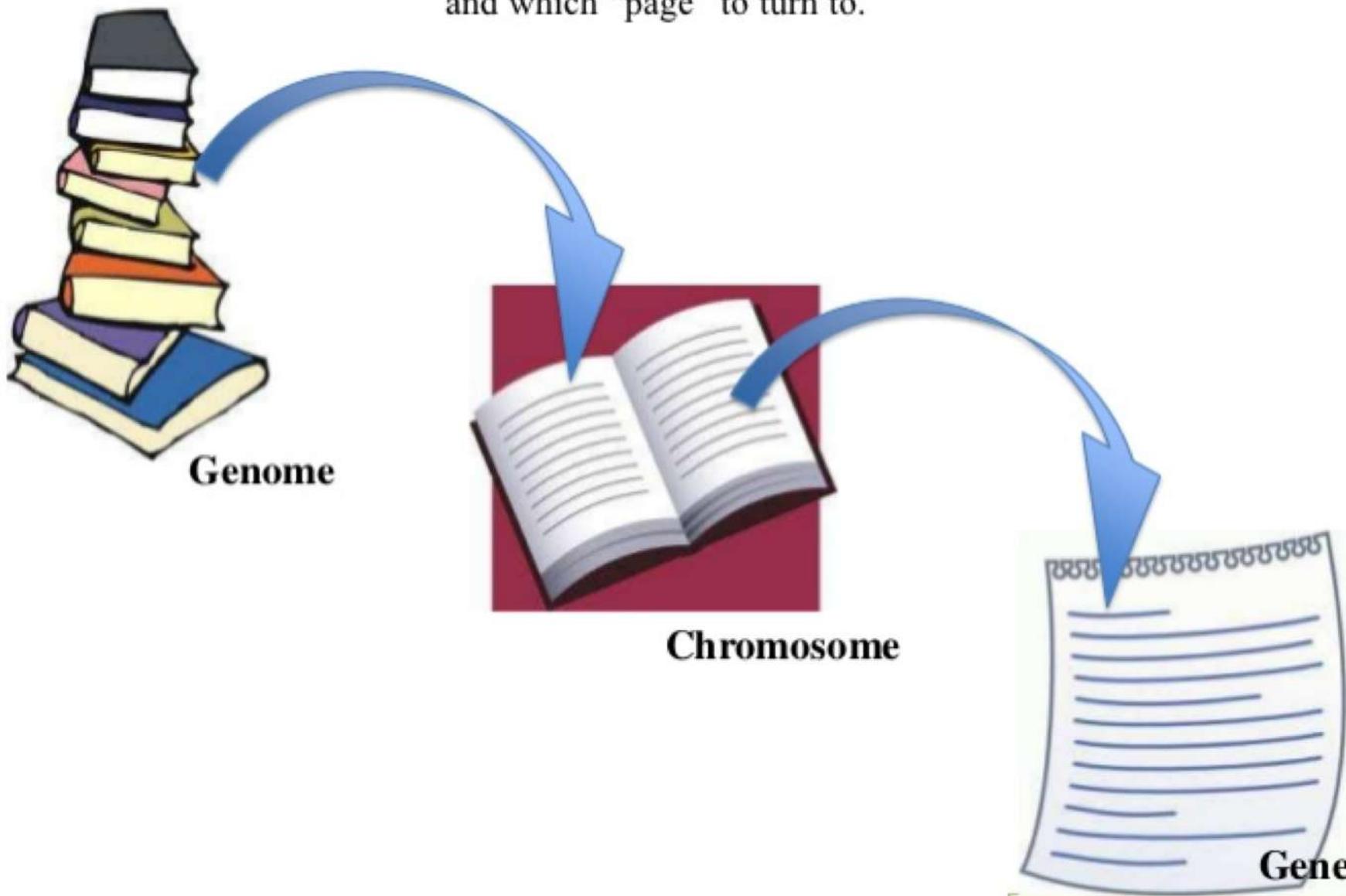


Our Genome:

ALL of the sets of instructions for how to make ALL the proteins we need

ALL (gene, chromosome, genome) are written in same the DNA alphabet!!!

In order for scientists to find a specific gene, they need to know which “book” to look in, and which “page” to turn to.



MONTAGEM DE UM GIGANTESCO QUEBRA-CABEÇA

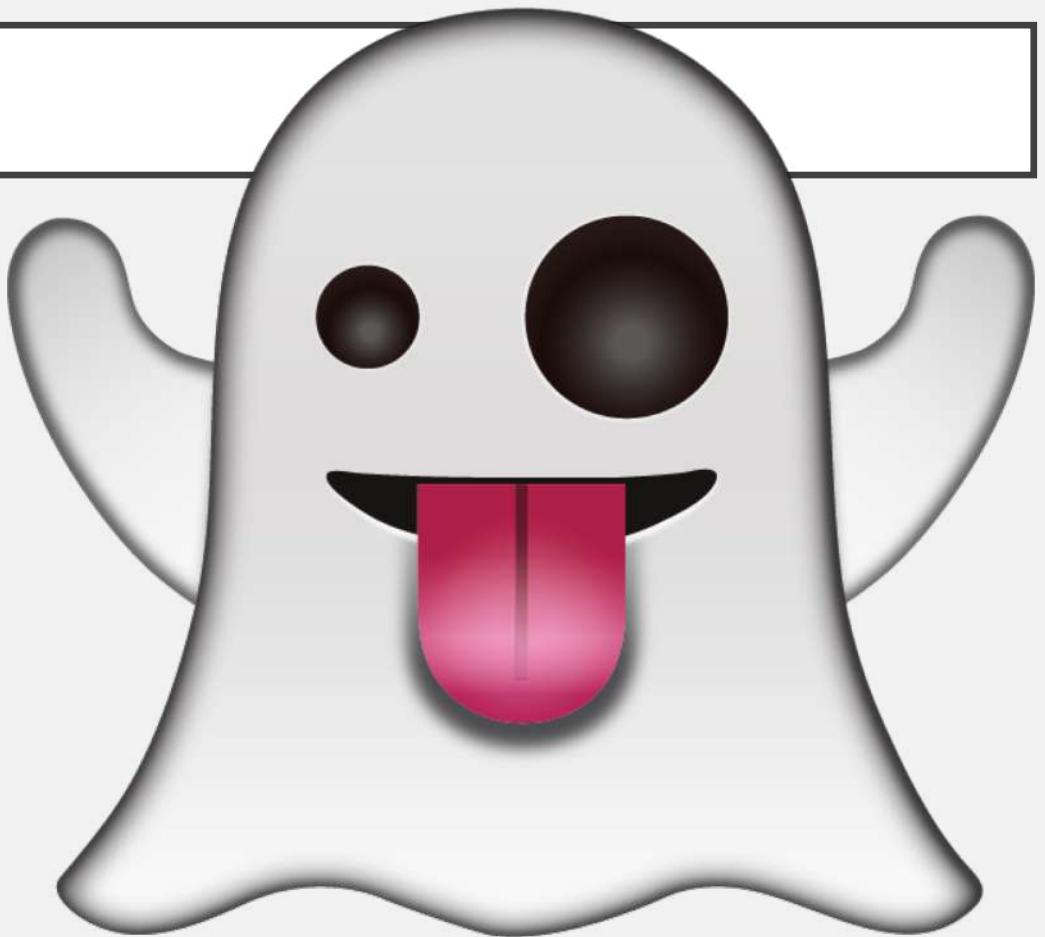
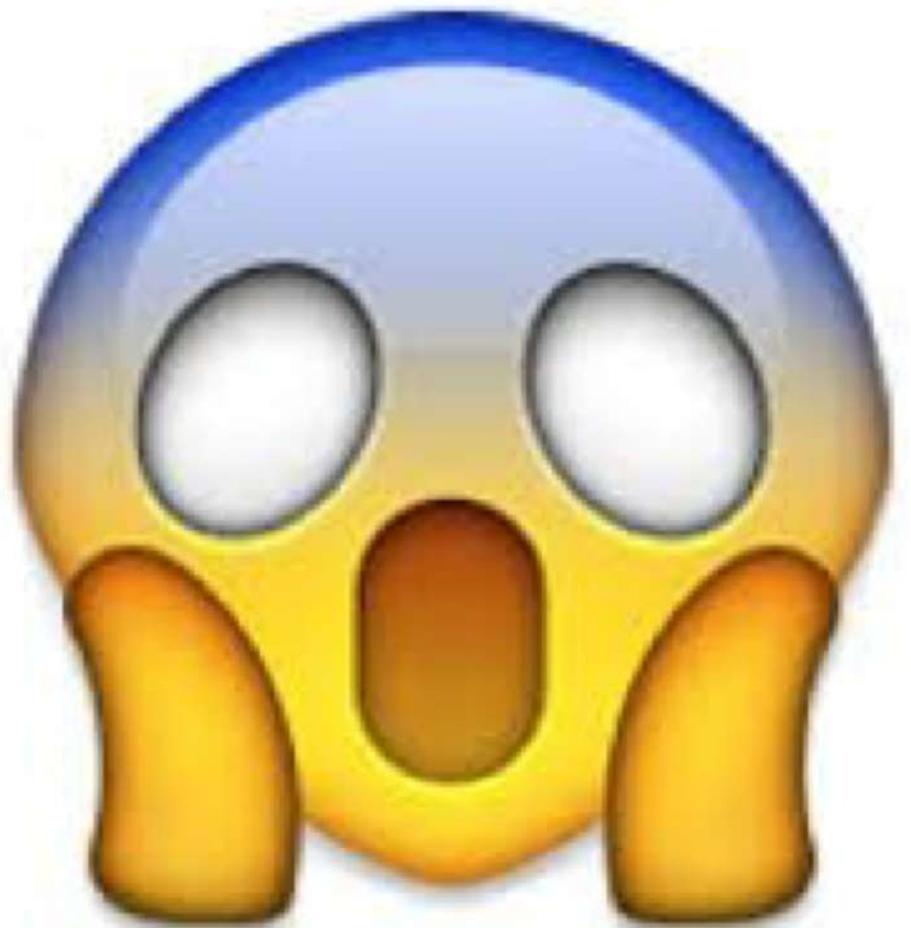


GIGANTESCO QUEBRA-CABEÇA SEM A CAIXA











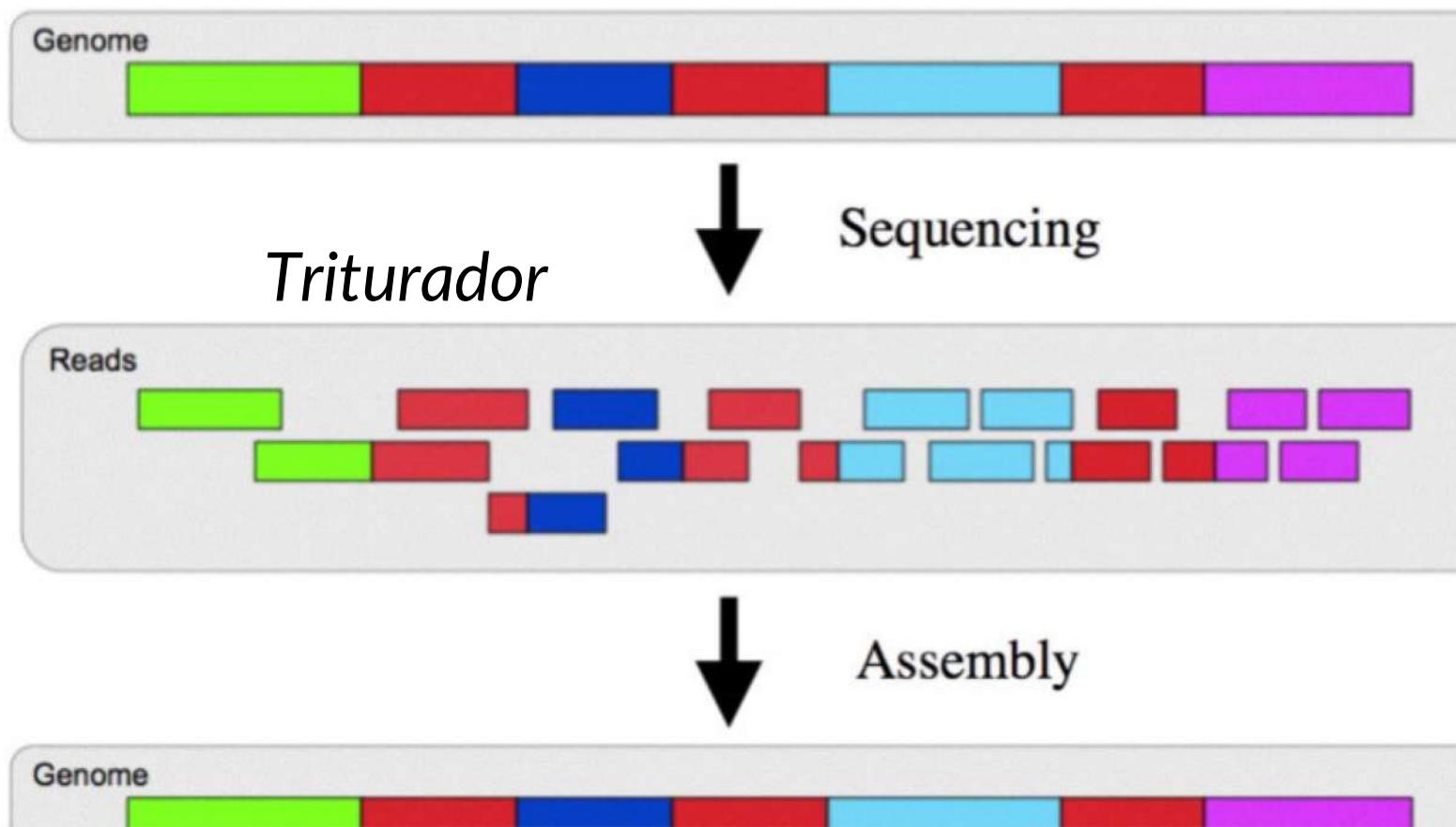
MONTAGEM DE NOVO

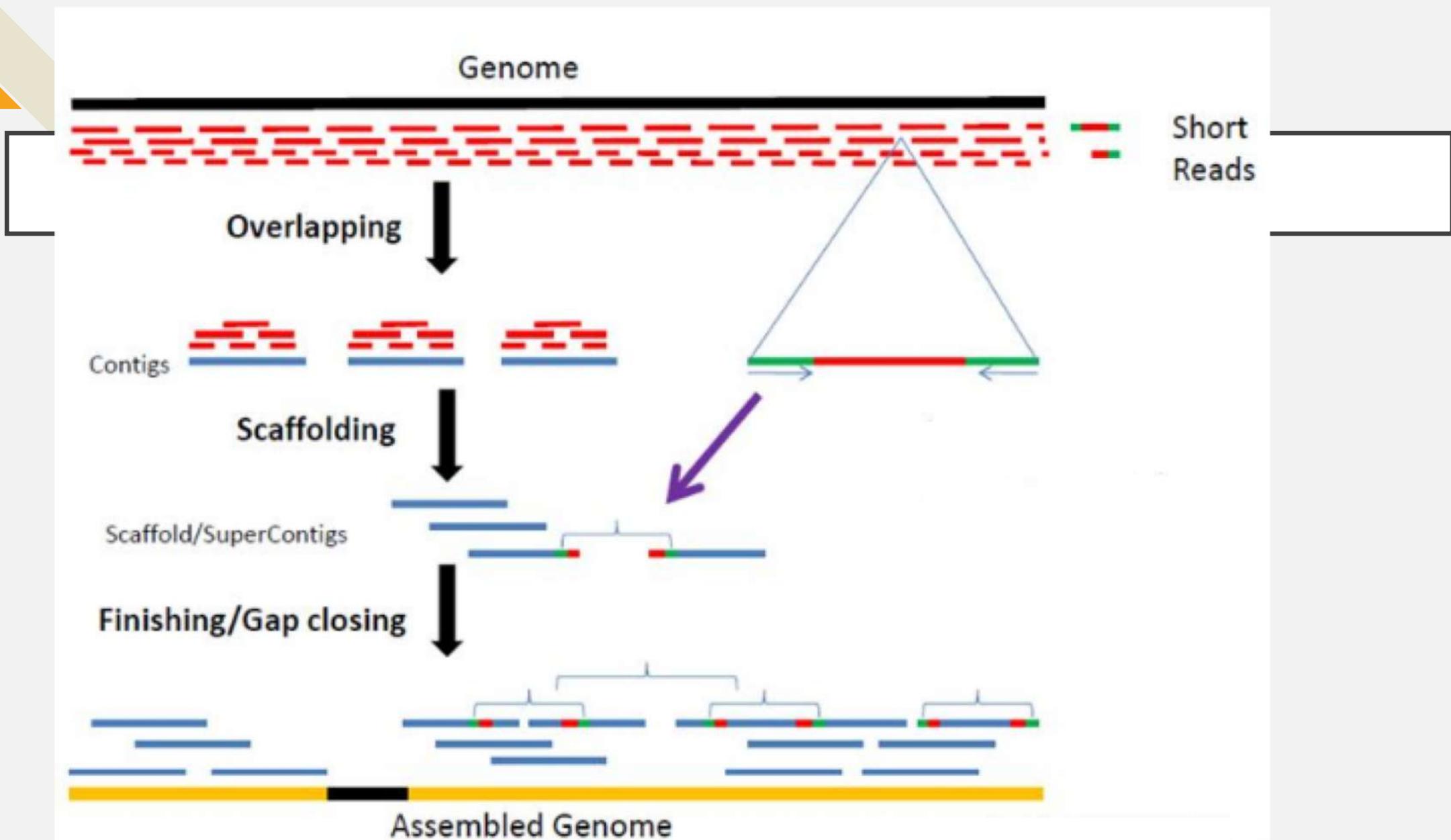
Método computacional para montar fragmentos de sequências de DNA
- usando sobreposição das extremidades

Leituras obtidas de métodos de sequenciamento para formar sequências de DNA mais longas

MONTAGEM DE GENOMA - BÁSICO

What is Genome Assembly?







MONTAGEM GENOMA REFERÊNCIA

Montagem usando sequência referência existente

Montagem apenas constrói uma sequência semelhante a referência

Algoritmo de mapeamento

Human Genome Sequencing

Generating a Reference Genome Sequence
(e.g., Human Genome Project)



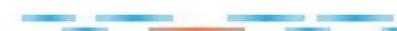
Generating a Person's Genome Sequence
(e.g., Circa ~2016)



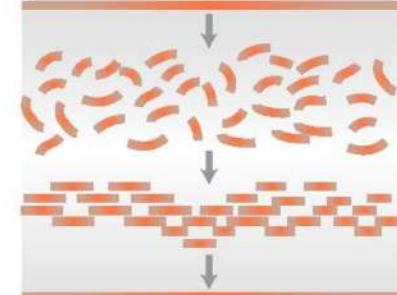
Break genome into large fragments and insert into clones



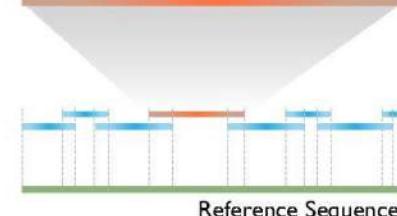
Order clones



Break individual clones into small pieces

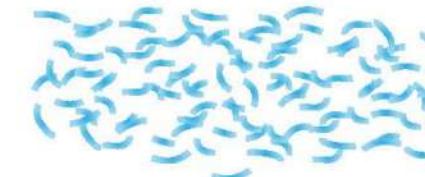


Generate thousands of sequence reads and assemble sequence of clone



Assemble sequences of overlapping clones to establish reference sequence

Genomic DNA



Break genome into small pieces

....TATGCGATGCGTATTCGTAAA....

Generate millions of sequence reads



Align sequence reads to established reference sequence

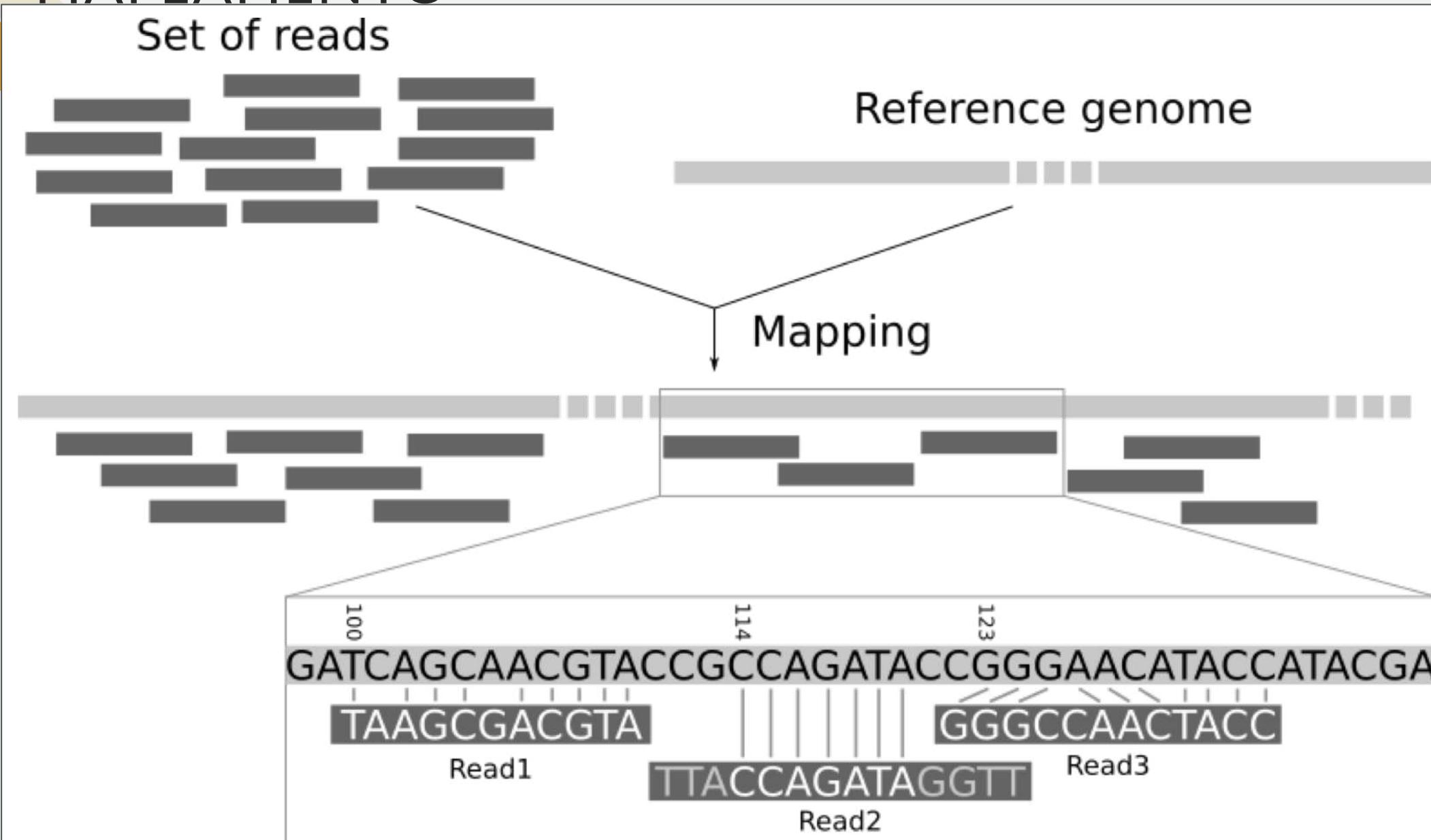
Reference Sequence

Deduce starting sequence and identify differences from reference sequence

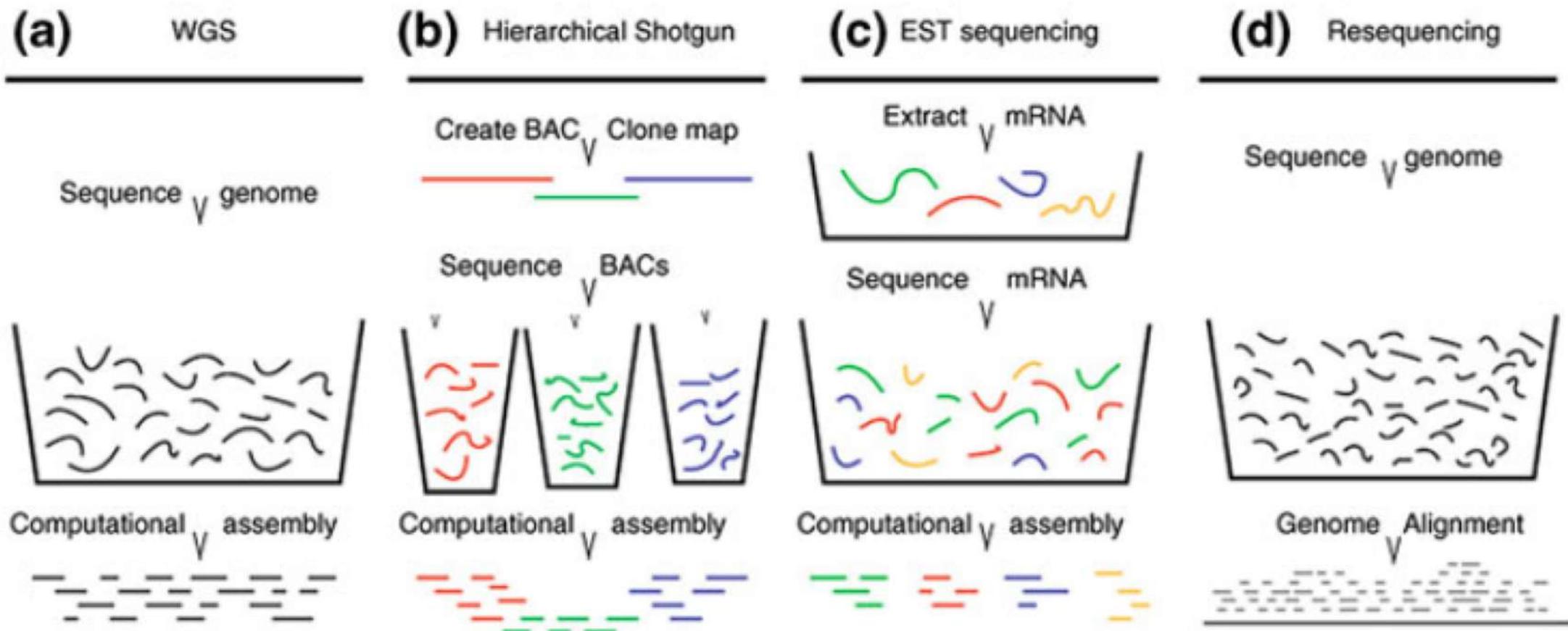
MONTAGEM GENOMA REFERÊNCIA



MONTAGEM DE GENOMA REFERÊNCIA - MAPEAMENTO



MÉTODOS DE SEQUENCIAMENTO E MONTAGEM



TAMANHO DE READS

TABLE I. Main sequencing technologies and their characteristics

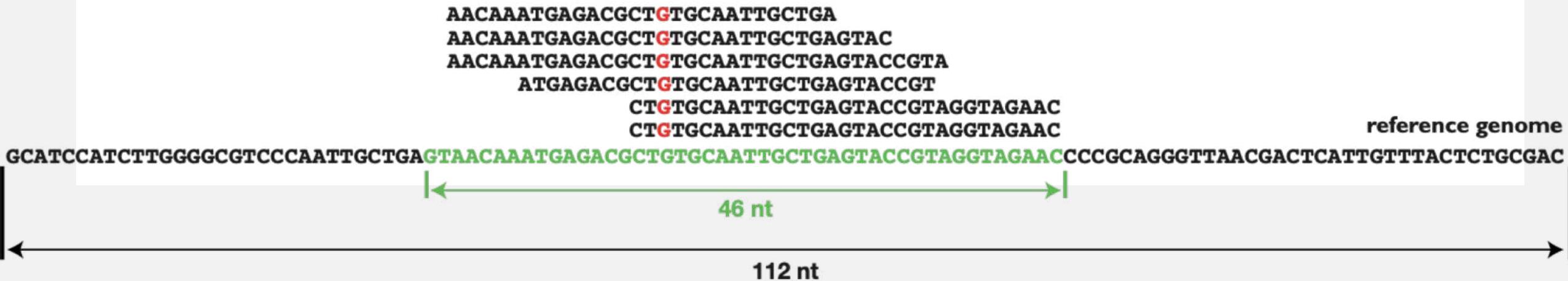
Technology (manufacturer)	Sequencing chemistry	Platform	Read length (bp)	Throughput (Gb/h run)	Best suited for:
Sanger	Dye terminator	ABI 3730xl	700–900		De novo and metagenomics
454 (Roche)	Pyrosequencing	GS FLX	400–700	0.04	De novo and metagenomics
		GS Junior	400	0.004	De novo and metagenomics
Solexa (Illumina)	Sequencing by synthesis with reversible terminators	Gaix HiSeq2000 MiSeq	36–150 36–100 36–250	0.3 2.9 0.2	Resequencing
SOLID (ABI)	Sequencing by ligation	5500xl	35–75	1	Resequencing
Heliscope (Helicos)	Sequencing by synthesis with virtual terminators	tSMS	25–55	1	Resequencing
Ion Torrent (Life Technologies)	Semiconductor sequencing	Ion torrent PGM	100–200	0.2	Resequencing
		Ion proton sequencer	100–200	2.5	Resequencing
PacBio (Pacific Bioscience)	SMRT technology	PacBioRS	250–10 000	0.1	Genome structure and metagenomics
Nanopore (Oxford Nanopore Technologies)	Ionic current sensing	GridION and MinION	10 000–50 000*	* *	De novo and genome structure

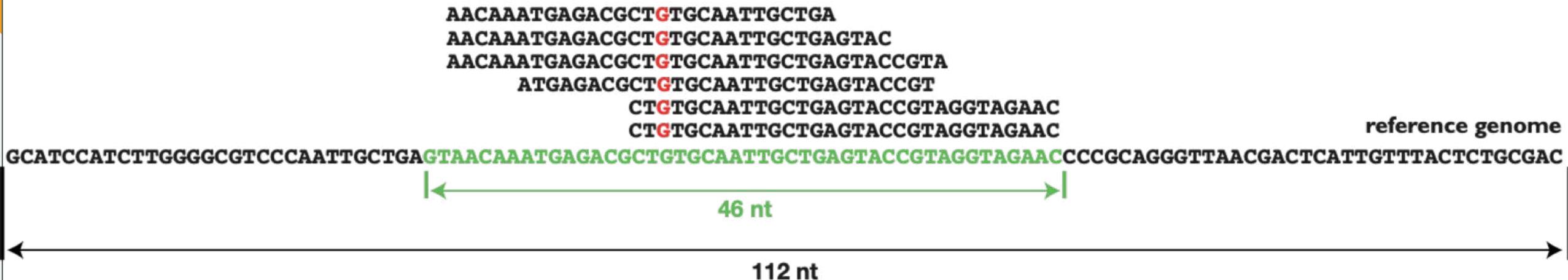
*Data not available on the corporate website.

Para minimizar o problema de tamanhos muito pequenos dos reads deve aumentar a cobertura

COBERTURA

Número de *reads* que representam um dado nucleotídeo na sequência do contig





Reference point	Calculation	Example (see Fig.2)
Whole genome	(# of sequenced bases) * / (genome size)	$188 / 112 = 1,68$ fold
One locus	(# of bases mapping to the locus) / (size of locus)	$188 / 46 = 4,09$ fold
One position	(# of reads overlapping with one position)	6 fold



K-MER (LEITURAS DE COMPRIMENTO K)

Estimativa de tendências do sequenciamento

Conteúdo de repetições



K-MER (LEITURAS DE COMPRIMENTO K)

Cobertura do sequenciamento

Detecta heterozigotos

Erros na montagem de genomas

Table 3.1 The mean number of false placements of k-mers on the genome [5]

K	<i>Escherichia coli</i>	<i>Saccharomyces cerevisiae</i>	<i>Arabidopsis thaliana</i>	<i>Homo sapiens</i>
200	0.063	0.26	0.053	0.18
160	0.068	0.31	0.064	0.49
120	0.074	0.39	0.086	1.7
80	0.082	0.49	0.15	7.2
60	0.088	0.58	0.27	18
50	0.091	0.63	0.39	32
40	0.095	0.69	0.65	78
30	0.11	0.77	1.5	330
20	0.15	1.0	5.7	2,100
10	18	63.8	880	40,000

Montagem incorreta do k-mers, depois de remover as montagens corretas



SOBREPOSIÇÃO DO READ

Reads que se sobrepõem a mais de um read formarão um ramo na montagem (incerteza)

Um dos métodos para abordar essas ramificações é usando leituras em pares (*paired-end reads*)

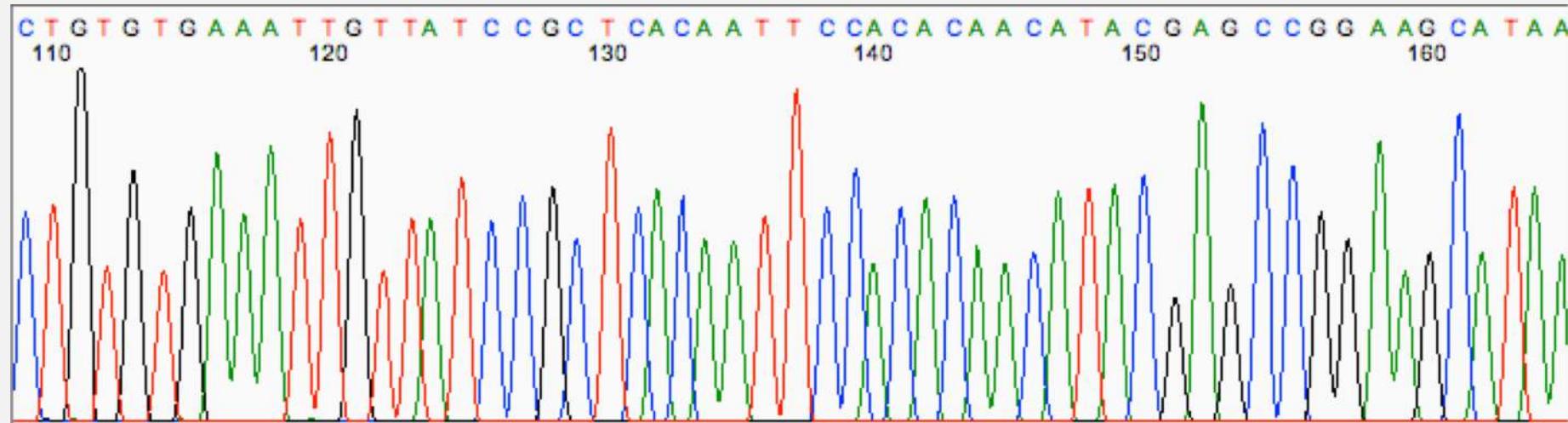


QUALIDADE DO READ

Os erros de sequenciamento podem atrapalhar o processo de montagem

O uso da informação da qualidade da base pode ajudar na montagem e descobrir erros de sequenciamento

QUALIDADE DA BASE NO READ



Nem todos os programas usam a qualidade das bases porque isso demanda mais esforço computacional



PROBLEMAS COM SEQUENCIAMENTO

Leituras (*reads*) curtas

Regiões de baixa cobertura



PROBLEMAS COM SEQUENCIAMENTO

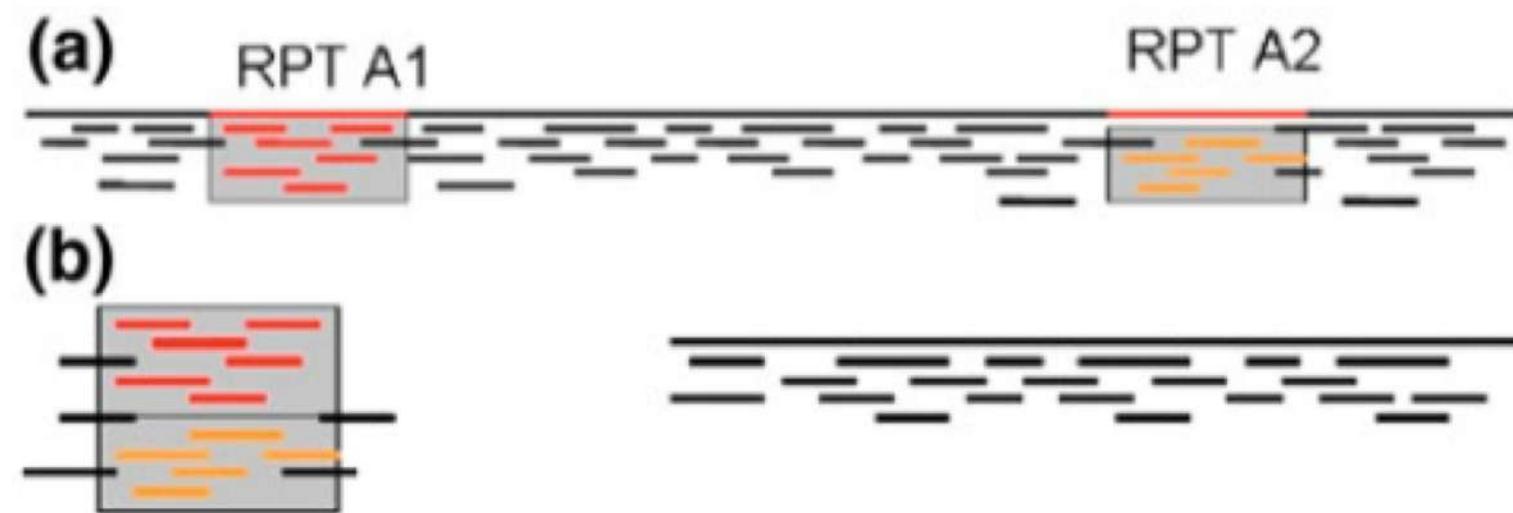
Regiões repetitivas

Erros de sequenciamento

PROBLEMAS COM SEQUENCIAMENTO

3.2 Challenges of Genome Assembly

45



Além disso, algumas partes do genoma podem permanecer não sequenciadas

REPETIÇÕES PODEM SER MINIMIZADAS COM ~~PAIR-END-READS~~

Pair-end-reads são fragmentos de DNA com ambas as extremidades sequenciadas.

Podem ser considerados como duas leituras com uma distância específica entre essas leituras.



PRÉ PROCESSAMENTO

Quando os dados não estão limpos e prontos para serem usados pelo algoritmo de montagem

Remover reads de baixa qualidade

Remover extremidades de baixa qualidade



COMO FILTRAR OS READS

Reads confiáveis : livres de erros e não-repetitivos

Se k-mer ocorre uma única vez , provavelmente é um erro de sequenciamento

Se k-mer ocorre muitas vezes , provavelmente é uma repetição

COMO FILTRAR OS READS

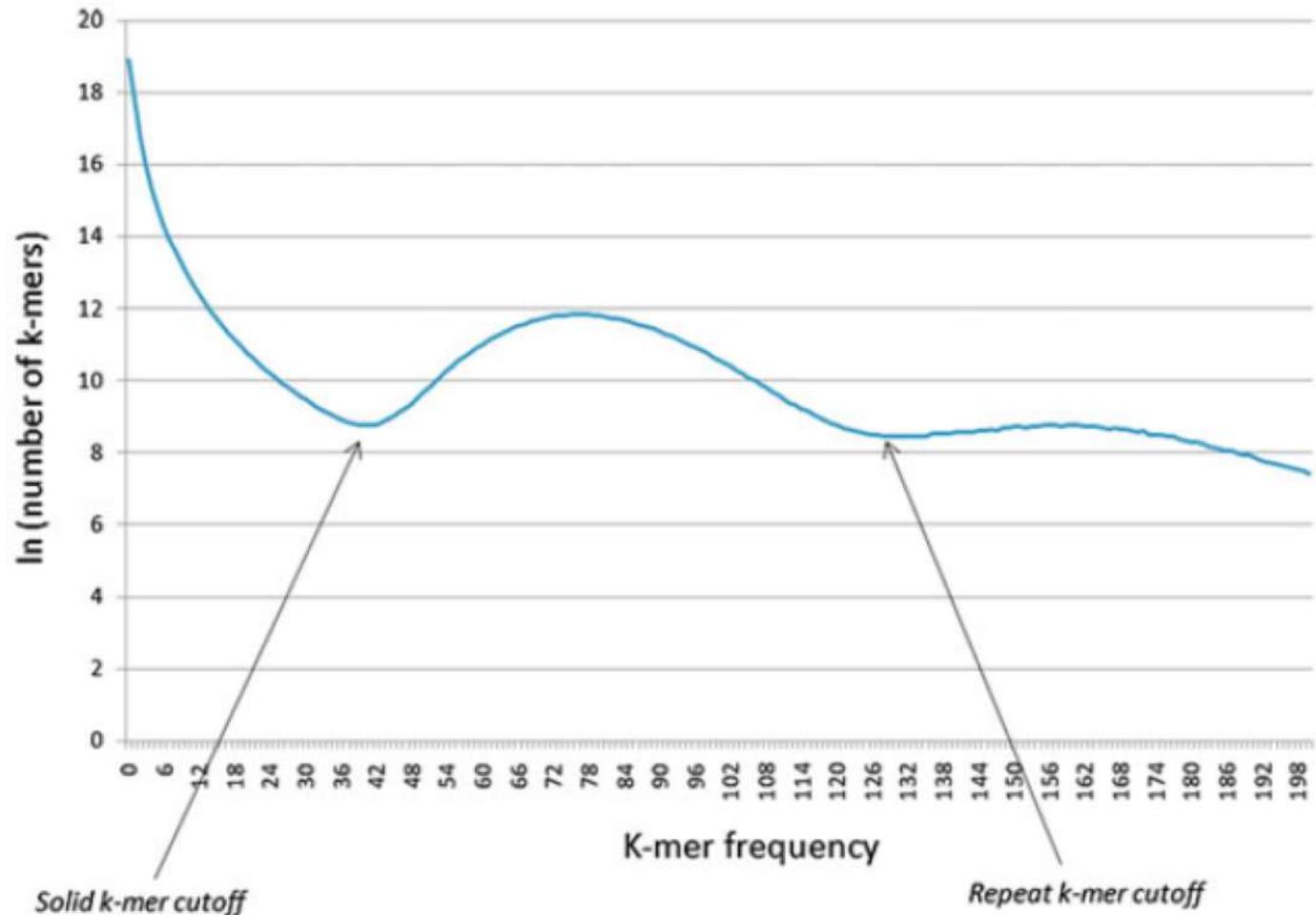


Fig. 3.4 Statistical analysis for filtering erroneous reads [6]

AVALIAÇÃO DA MONTAGEM

Pontuação de acerto com base em um genoma referência

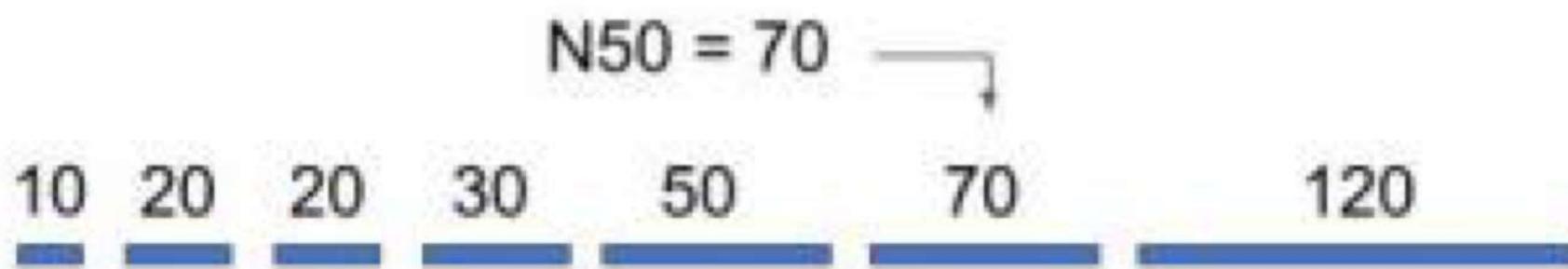
Estatísticas de tamanho, N50

DISTRIBUIÇÃO DO COMPRIMENTO CONTIG E SCAFFOLD

N50: estatística para análise de comprimentos de contig

N50: qual o menor comprimento de contigs que cobrem pelo menos metade do genoma

N50



Total length $10+20+20+30+50+70+120 = 320$

$50\% = 160$

$10+20+20+30+50+70 = 200$



ALGORITMO DE MONTAGEM

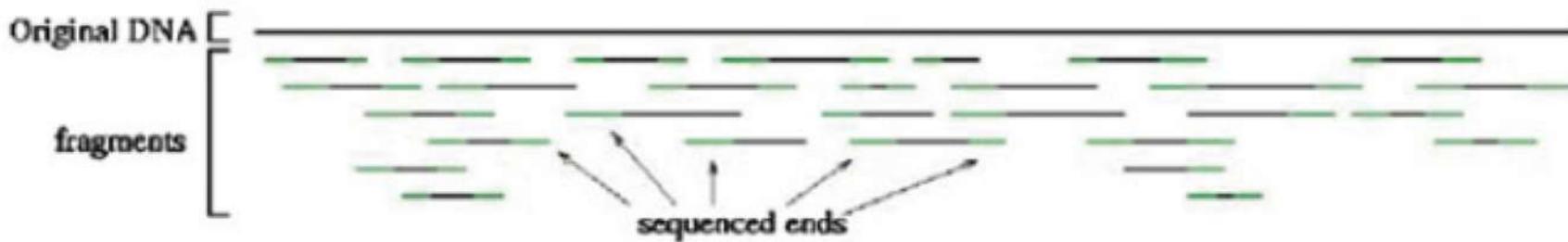
Algoritmo de montagem deve sobrepor os reads formando uma sequência longa e contínua (cromossomo)

Devido a erros de sequenciamento e à reads não sequenciados, contigs obtidos não são completos o suficiente e não formam um cromossomo.

Passos no processo de montagem

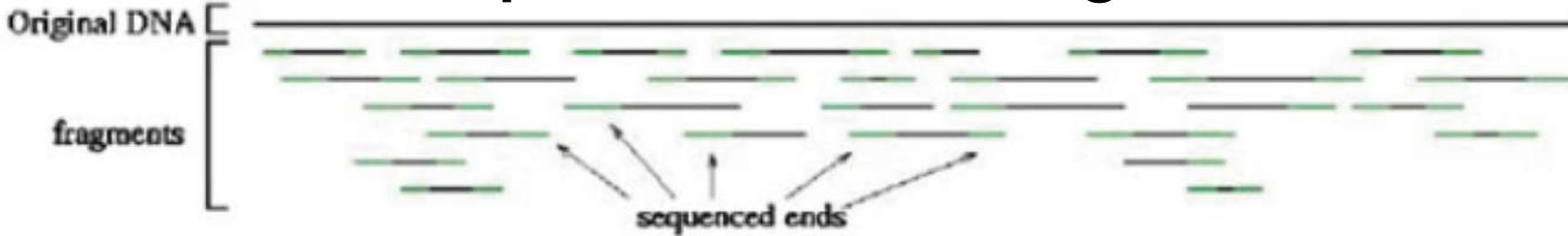


(a) Original DNA broken into a collection of fragments

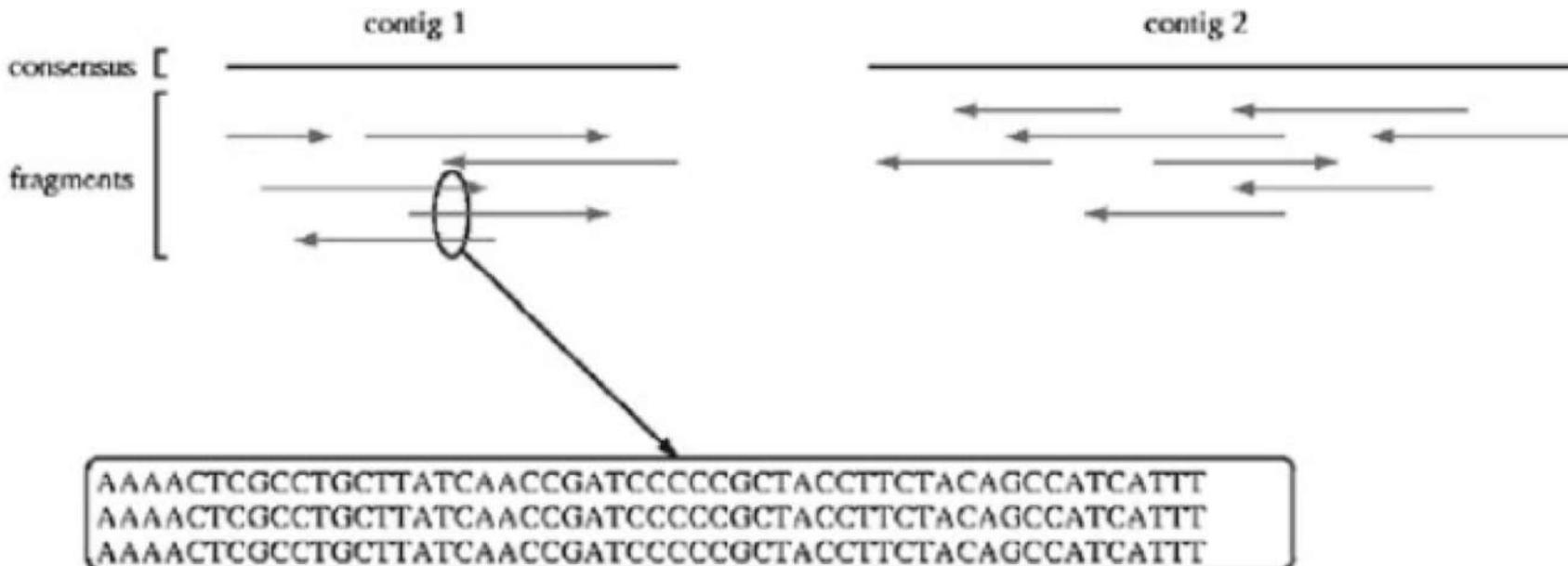


(b) The end of each fragment (drawn in green) are sequenced

Passos no processo de montagem



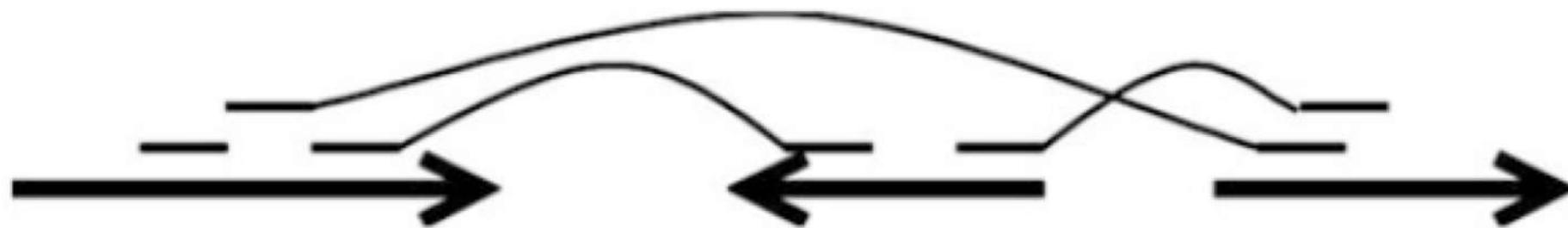
(b) The end of each fragment (drawn in green) are sequenced



(c) Merging reads to form contigs.

Passos no processo de montagem

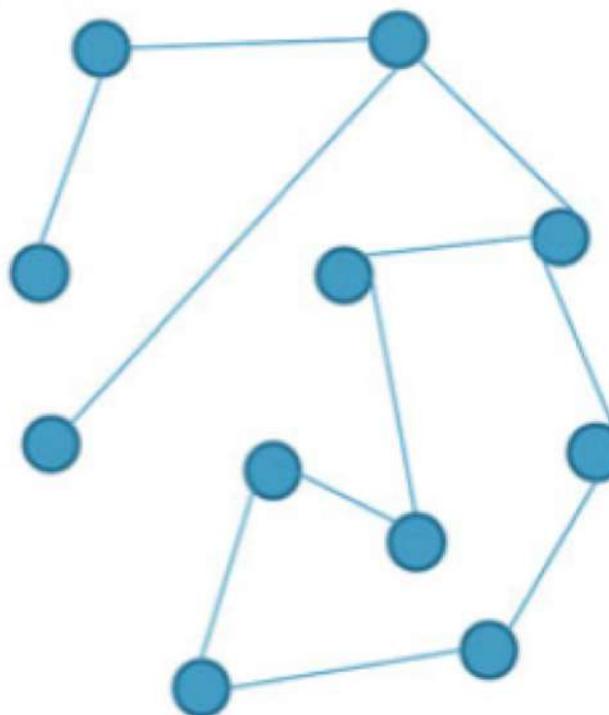
PAIR-END-READS



Necessário existir 3-5 pair-end-read com distância esperada para suportar um scafflod

MONTAGEM - PROBLEMA DOS GRAFOS

ATTGCC**CGGAA** **CGGAA**TGTGAT

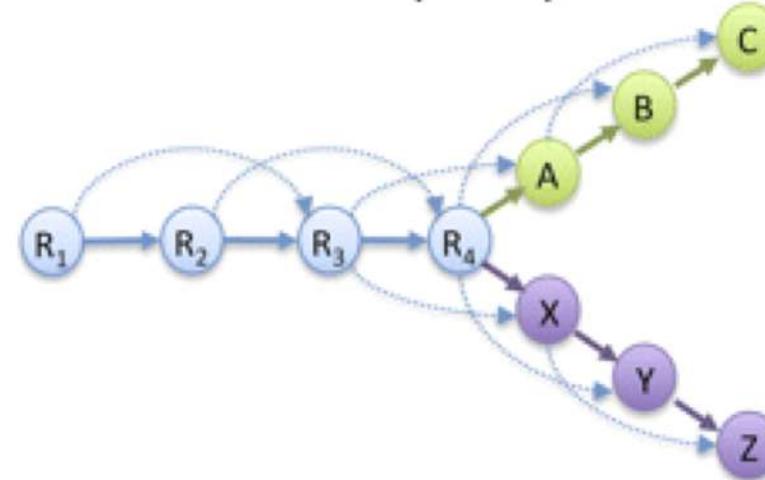


MONTAGEM SOBREPOSIÇÃO E GRAFO BRUIJN

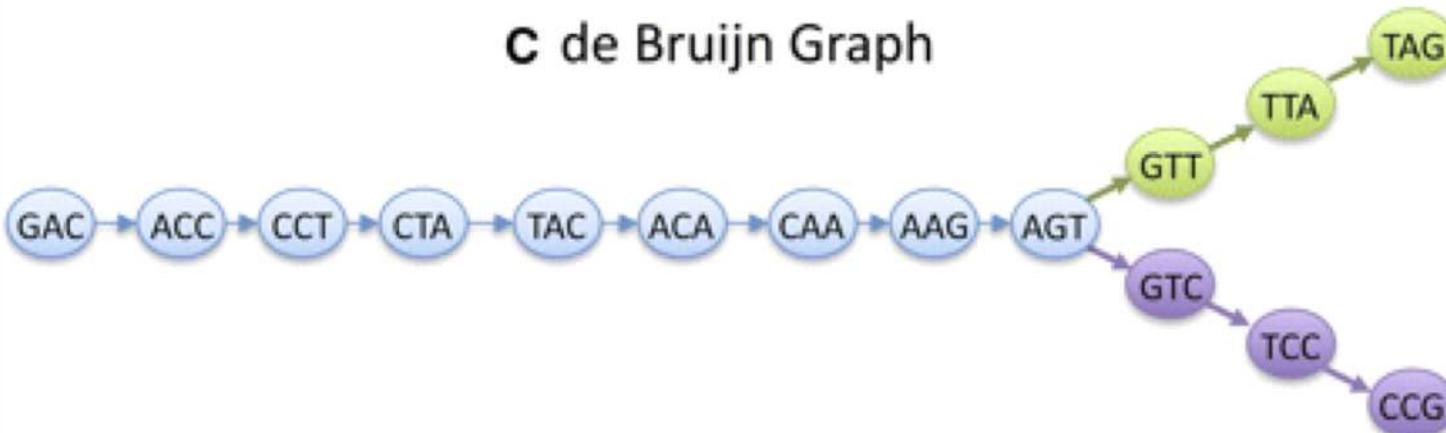
A Read Layout

R ₁ :	GACCTACA
R ₂ :	ACCTACAA
R ₃ :	CCTACAAG
R ₄ :	CTACAAAGT
A:	TACAAGTT
B:	ACAAGTTA
C:	CAAGTTAG
X:	TACAAGTC
Y:	ACAAGTCC
Z:	CAAGTCCG

B Overlap Graph



C de Bruijn Graph



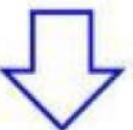
VELVET

Velvet é um algoritmo para montagem *de novo*

Adequado para montagem de *reads* curtos (25–50 bp)

N50 de 50kb para genomas procariotos

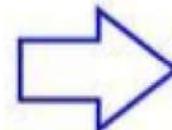
A A T G C C G T A C G T A G G G G T A A T A T A T G A C C A



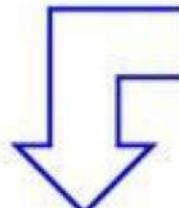
(Sequencing: Solexa, Illumina, etc..)

T G C C G T T A G G G T A T A T A T
A A T G C T T A C G T A A T G A C C
T T G C C G C G T A G G T A A T A T
G T A C G T G T A C T A
A A T G C C G G G T A A T G A C C A
 G T A G G G T A T G A C
 C T A T A T

Compute k-mer
with $k=4$

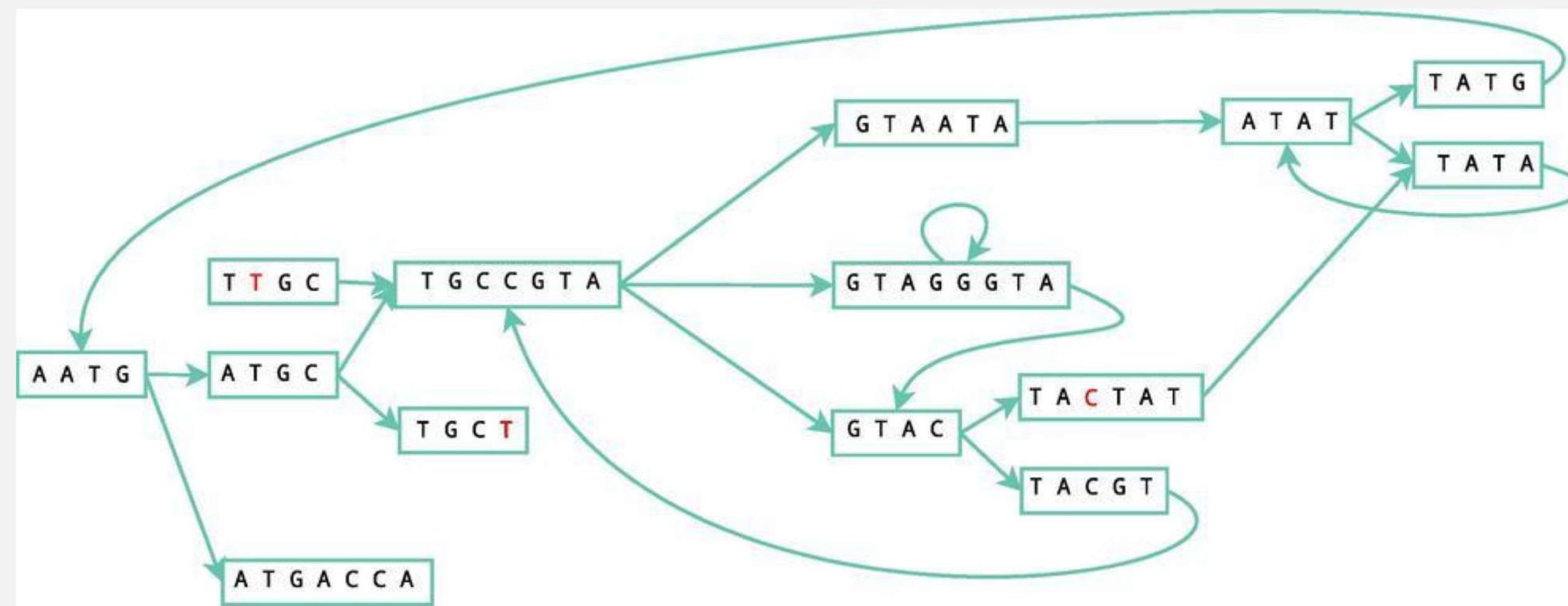


Create Graph for
the set of k-mers



A A T A
A A T G (x2)
A C C A
A C G T (x2)
A C T A
A G G G (x2)
A T A T (x4)
A T G A (x2)
A T G C (x2)
C C G T
C G T A (x2)
C T A T
G A C C (x2)
G C C G (x2)
G G G T (x2)
G G T A
G T A A
G T A C (x2)
G T A G (x2)
T A A T
T A C G (x2)
T A C T
T A G G (x3)
T A T A (x2)
T A T G
T G A C (x3)
T G C C (x3)
T G C T
T T G C

VELVET



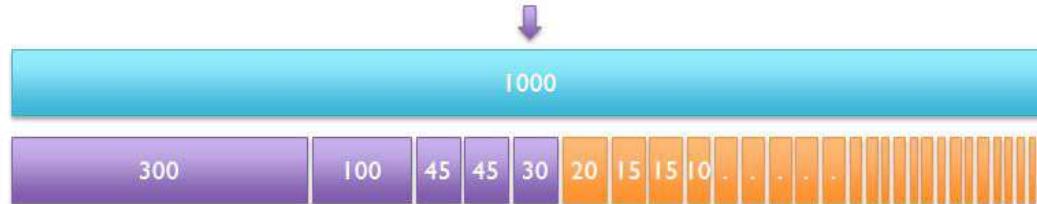
N50

N50 size

Def: 50% of the genome is in contigs as large as the N50 value

Example: 1 Mbp genome

50%



N50 size = 30 kbp

$$(300k + 100k + 45k + 45k + 30k = 520k \geq 500\text{ kbp})$$

A greater N50 is indicative of improvement in every dimension:

- Better resolution of genes and flanking regulatory regions
- Better resolution of transposons and other complex sequences
- Better resolution of chromosome organization
- Better sequence for all downstream analysis