

SCIENCES

BIOLOGY

Genetics, Epigenetics

Function and Evolution of Repeated DNA Sequences

**Coordinated by
Guy-Franck Richard**

ISTE

WILEY

Table of Contents

[Cover](#)

[Table of Contents](#)

[Title Page](#)

[Copyright Page](#)

[Foreword](#)

[Introduction: About Repeated Genomes](#)

[I.1. The “C-value” paradox](#)

[I.2. Recycling junk DNA](#)

[I.3. The different repeat types](#)

[I.4. References](#)

[1 Whole-Genome Duplications, a Source of Redundancy at the Entire-Genome Scale](#)

[1.1. Prevalence of polyploids in the tree of life](#)

[1.2. Mechanisms for the appearance of whole-genome duplications](#)

[1.3. Cellular consequences of whole-genome duplications](#)

[1.4. Rediploidization: evolutionary reduction in genetic redundancy](#)

[1.5. Functions and evolution of duplicated genes](#)

[1.6. Whole-genome duplications and evolutionary diversification](#)

[1.7. Perspectives and conclusions](#)

[1.8. References](#)

[2 Segmental Duplications and CNVs: Adaptive Potential of Structural Polymorphism](#)

[2.1. The multiple facets of genetic polymorphism](#)

[2.2. From Segmental Duplications to Copy Number Variants: terminology](#)

[2.3. SDs: a general overview](#)

[2.4. Methodologies for detecting structural variation in genomes](#)

[2.5. The molecular mechanisms at the origin of structural variation](#)

[2.6. Regions rich in SDs/LCRs favor the creation of CNVs: insertions/ duplications, deletions and inversions](#)

[2.7. From SDs to CNVs in humans and primates](#)

[2.8. SDs in little-studied species: general genomic profiles](#)

[2.9. SD content: impact of a duplicated environment on sequences that make up the SDs](#)

[2.10. SDs and epigenetic modifications](#)

[2.11. The adaptive potential of SDs: between the benefit of innovation and the cost of pathology](#)

[2.12. SDs and associated CNVs: their roles in species adaptation to changes in environments](#)

[2.13. Conclusion](#)

[2.14. Glossary of terms](#)

[2.15. References](#)

[3 Transposable Elements: Parasites that Shape Genome Evolution](#)

[3.1. Transposable elements in eukaryotic genomes](#)

[3.2. Classification of TEs and transposition mechanisms](#)

[3.3. TE self-regulation](#)

[3.4. TE restriction by the host](#)

[3.5. The impact of transposition events on genomes](#)

[3.6. Conclusion](#)

[3.7. References](#)

[4 Insights Into the Evolutionary Diversity of Centromeres](#)

[4.1. The centromere](#)

[4.2. Monocentromeres](#)

[4.3. Holocentromeres](#)

[4.4. Open questions](#)

[4.5. Acknowledgments](#)

[4.6. References](#)

[5 Evolution and Functions of Telomeres](#)

[5.1. Primary structure of telomeres](#)

[5.2. A telomere specific higher order structure: the T-loop](#)

[5.3. Telomere lengthening mechanisms](#)

[5.4. Telomere length homeostasis](#)

[5.5. Telomeres and genome organization and function](#)

[5.6. Cell senescence, aging and disease](#)

[5.7. Conclusion](#)

[5.8. Acknowledgments](#)

[5.9. References](#)

[6 G-quadruplexes: Structure, Detection and Functions](#)

[6.1. From guanine-guanine base-pairing to a secondary structure](#)

[6.2. The G4 structure: variations on a theme](#)

[6.3. Finding G-quadruplexes in a genome](#)

[6.4. Biological roles of G-quadruplexes](#)

[6.5. Perspective: G-quadruplexes as anticancer therapeutic targets](#)

[6.6. References](#)

[7 Satellite DNA, Microsatellites and Minisatellites](#)

[7.1. Satellite DNAs, origin and definition](#)

[7.2. From semantics to biology](#)

[7.3. The evolutionary mechanisms of tandem repeats](#)

[7.4. Microsatellites in human diseases](#)

[7.5. De novo formation and evolution of tandem repeats](#)

[7.6. Perspectives](#)

[7.7. Acknowledgments](#)

[7.8. References](#)

[8 CRISPR-Cas: An Adaptive Immune System](#)

[8.1. A brief history of the discovery of CRISPR-Cas systems](#)

[8.2. General characteristics of CRISPR-Cas systems](#)

[8.3. Evolution of CRISPR-Cas systems](#)

[8.4. An adaptive immune system](#)

[8.5. Phage escape mechanisms](#)

[8.6. Biological cost of CRISPR-Cas systems](#)

[8.7. Importance in nature: impact of ecological factors](#)

[8.8. Conclusions and perspectives](#)

[8.9. References](#)

[List of Authors](#)

[Index](#)

[End User License Agreement](#)

List of Tables

Chapter 2

[Table 2.1. Summary of SD detection methods](#)

[Table 2.2. Information on the SDs extracted from 12 genomes](#)

[Table 2.3. SDs in humans - Gene, modification, phenotype-disease, mechanism an...](#)

[Table 2.4. CNVs and examples of traits affected during domestication](#)

Chapter 3

[Table 3.1. Examples of “transposopathies” for which TE insertion is associated...](#)

[Table 3.2. Some notable or iconic examples of domestication of different TE en...](#)

[Table 3.3. Examples of natural transposition events, selected by humans](#)

[Table 3.4. Examples of molecular biology tools developed from TE](#)

Chapter 6

[Table 6.1. Biophysical and biochemical methods to study G-quadruplexes \(for ad...](#)

[Table 6.2. Putative G-quadruplex sequences in 12 genomes \(for a recent and det...](#)

[Table 6.3. Resolved G-quadruplex structures in gene promoters](#)

List of Illustrations

Introduction

[Figure I.1. Example of \$C_0t\$ curve](#)

[Figure I.2. Comparison of genome sizes and gene numbers](#)

[Figure I.3. The different types of repeated DNA sequences](#)

Chapter 1

[Figure 1.1. Whole-genome duplications identified in the eukaryotic phylogeneti...](#)

[Figure 1.2. Tetraploidization by endoreplication. In the case of a normal cell...](#)

[Figure 1.3. Allopolyploidizations in the lineage of wheat \(*Triticum aestivum*\)....](#)

[Figure 1.4. Restoration of meiosis in polyploids. In polyploid species, homolo...](#)

[Figure 1.5. Organization of homeologous regions in the genome of rainbow trout...](#)

[Figure 1.6. Homeologous regions with double-conserved synteny in teleost fish ...](#)

[Figure 1.7. Ancestral and delayed rediploidization. In the most classical case...](#)

[Figure 1.8. Outcome of duplicated genes after a polyploidization event. During...](#)

[Figure 1.9. Example of ohnologous genes with divergent territories of expressi...](#)

[Figure 1.10. Preferential retention mechanisms of ohnolog copies of genes. The...](#)

[Figure 1.11. Disentangling of a gene regulatory block during rediploidization....](#)

Chapter 2

[Figure 2.1. General view of mutations affecting eukaryotic genomes, their impa...](#)

[Figure 2.2. The human Y chromosome, its major duplications, and the alteration...](#)

[Figure 2.3. Representation of interchromosomal \(center\) and intrachromosomal \(...\)](#)

[Figure 2.4. Distribution of duplications on human Y chromosomes compared to a ...](#)

[Figure 2.5. Usual fates of a duplicated gene: \(i\) conservation; \(ii\) subfuncti...](#)

[Figure 2.6. Rearrangements involving SDs that are directly \(duplication/deleti...](#)

Chapter 3

[Figure 3.1. Proportion of TE sequences in several genomes. While the proportio...](#)

[Figure 3.2. Classification of autonomous TEs in eukaryotes. This classificatio...](#)

[Figure 3.3. Main mechanisms controlling the expression and mobility of TEs. Ev...](#)

[Figure 3.4. Examples of integration sites targeted by TEs. Chromosome features...](#)

[Figure 3.5. Recognition and repression of TEs by the KRAB-ZFP complex. The KRA...](#)

[Figure 3.6. Biogenesis of piRNAs derived from uni-strand piRNA clusters. Uni-s...](#)

[Figure 3.7. Biogenesis of piRNAs derived from dual-strand piRNA clusters in Dr...](#)

[Figure 3.8. Mutagenesis induced by the insertion of a TE in a gene. Shown at t...](#)

[Figure 3.9. Functional consequences of the insertion of a TE in the regulatory...](#)

[Figure 3.10. Consequences of recombination between TEs of the same family \(NAH...](#)

[Figure 3.11. Complex interactions between TEs and their host. TE activity is r...](#)

Chapter 4

[Figure 4.1. \(A\) Illustrations of salamander cells with chromosomes \(monocentri...](#)

[Figure 4.2. Phylogeny of several fungal organisms along with their respective ...](#)

[Figure 4.3. Schematics of models for different holocentric architectures of C....](#)

Chapter 5

[Figure 5.1. On top, canonical nucleotide sequence of vertebrate telomeres. The...](#)

[Figure 5.2. Emergence and evolution of linear chromosomes \(simplified from Vil...](#)

[Figure 5.3. The nucleoprotein structure of human telomeres. The telomere-speci...](#)

[Figure 5.4. Mechanisms of telomere replication](#)

[Figure 5.5. The end-replication problem at telomeres is a leading mechanism pr...](#)

[Figure 5.6. Telomere length shortens with age and exaggerated shortening is as...](#)

Chapter 6

[Figure 6.1. From guanines to the G-quadruplex structure. From left to right: a...](#)

[Figure 6.2. Strand orientation and G-quadruplex topologies. The glycosidic bon...](#)

[Figure 6.3. Parallel-stranded DNA G-quadruplex model reconstruction. The UCSF ...](#)

[Figure 6.4. Examples of loop conformations from published G-quadruplex structu...](#)

[Figure 6.5. Targeting G-quadruplexes. The model reconstruction of a G-quadrupl...](#)

Chapter 7

[Figure 7.1. Length distribution of the different tandem repeats. The abscissa ...](#)

[Figure 7.2. Number of citations in the PubMed database for different tandem re...](#)

[Figure 7.3. Microsatellite expansion disorders. On the left: diseases are clas...](#)

[Figure 7.4. Replication slippage model. During DNA synthesis in a tandemly rep...](#)

[Figure 7.5. Model of slippage during homologous recombination](#)

[Figure 7.6. Mechanisms of microsatellite formation by mutation. An unrepeated ...](#)

[Figure 7.7. Mechanisms of microsatellite formation by end-joining. A double-st...](#)

[Figure 7.8. Mechanism of formation of a minisatellite. A slippage between two ...](#)

[Figure 7.9. Mechanisms of evolution of mini- and megasatellites. \(A\) The dupli...](#)

[Figure 7.10. The limits of microsatellite detection algorithms. Depending on t...](#)

Chapter 8

[Figure 8.1. \(A\) Representation of the CRISPR genomic sequence with its success...](#)

[Figure 8.2. Evolution of the CRISPR array following the infection of a sensiti...](#)

[Figure 8.3. Conservation of the secondary structure of repeats of the same gro...](#)

[Figure 8.4. \(A\) Genetic organization of class 1 and class 2 systems. \(B\) Modul...](#)

[Figure 8.5. The three stages of the immune response. Adaptation corresponds to...](#)

[Figure 8.6. Diversity of molecular mechanisms of type I, III and II systems du...](#)

SCIENCES

Biology, Field Director – Marie-Christine Maurel

Genetics, Epigenetics, Subject Head – Bernard Dujon

Function and Evolution of Repeated DNA Sequences

Coordinated by

Guy-Franck Richard

ISTE

WILEY

First published 2023 in Great Britain and the United States by ISTE Ltd and John Wiley & Sons, Inc.

Apart from any fair dealing for the purposes of research or private study, or criticism or review, as permitted under the Copyright, Designs and Patents Act 1988, this publication may only be reproduced, stored or transmitted, in any form or by any means, with the prior permission in writing of the publishers, or in the case of reprographic reproduction in accordance with the terms and licenses issued by the CLA. Enquiries concerning reproduction outside these terms should be sent to the publishers at the undermentioned address:

ISTE Ltd
27-37 St George's Road
London SW19 4EU
UK

www.iste.co.uk

John Wiley & Sons, Inc.
111 River Street
Hoboken, NJ 07030
USA

www.wiley.com

© ISTE Ltd 2023

The rights of Guy-Franck Richard to be identified as the author of this work have been asserted by him in accordance with the Copyright, Designs and Patents Act 1988.

Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s), contributor(s) or editor(s) and do not necessarily reflect the views of ISTE Group.

Library of Congress Control Number: 2022949377

British Library Cataloguing-in-Publication Data
A CIP record for this book is available from the British Library
ISBN 978-1-78945-119-1

ERC code:
LS2 Genetics, 'Omics', Bioinformatics and Systems Biology

LS2_5 Epigenetics and gene regulation

Foreword

Our modern societies are too preoccupied with immediate performances to conceive of a world where the costs and efforts to achieve a result are not rationally minimized. And yet, life offers this image as soon as we take time to study it closely. The remarkable adaptation of different organisms to their living conditions revolves around genomes which are far from the products of what we may consider to be rational engineering. There is no such thing as a minimal genome: all of them are too large in comparison to the number of genes considered necessary to produce the organism hosting them. Often far too large, ours appears 50 times too large. They all contain identically (or almost identically) repeated sequences, sometimes they are repeated numerous times in the same genome; whereas random combinations of the four nucleotides make this phenomenon extremely unlikely, if not practically impossible.

This situation was observed as far back as the mid-20th century, long before the emergence of genomics, through the study of the renaturation kinetics of DNA molecules. The excessive amount of DNA and the abundance of repeated sequences remained a puzzle that some tended to quickly dismiss by referring to junk DNA, continuing to focus their studies only on what they already knew! Genome sequencing would come to solve this enigma by demonstrating just how incomplete our prior knowledge was. A new vision of genome organization and function is now provided, in which temporal dynamics combine with the present, because all genomes are simply imperfect copies of the genomes that preceded them and not new constructs. From now on, traces of the past mix with

present events and together they lay the foundations of the future.

As the present work illustrates so remarkably, repeats found in genomes can result from major evolutionary accidents such as whole-genome duplications, which, in a singular phenomenon, tend to coincide with the transitions between major geological eras. But they may also come from repeated interactions with infectious elements - of the viral type - that eventually integrate into chromosomes and are transmitted to the offspring. Thus, there are both endogenous and exogenous causes for the existence of repeated sequences in genomes. In turn, repeats can form the basis of the formation of chromosome functional elements such as centromeres, telomeres and guanine quadruplexes. Copy number variation of long repeated sequences can play a critical role in phenotypes and in organism adaptation. Similarly, the instability of short-sequence repeats allows us to easily differentiate between individuals from the same population. However, this can sometimes lead to very serious syndromes. Finally, the different mobile elements, kinds of specialized molecular machines, present in various numbers in the different genomes cannot be ignored. By the mid-20th century, they had already been identified by their genetic effects - they are mutagenic - but we now have a much broader view of their diversity and of the consequences, sometimes considerable, of their activity.

If our knowledge of repeated genome sequences has only progressed belatedly, this is partly due to technical difficulties encountered in their sequencing. Until the recent advent of new technologies allowing longer reads, it was very difficult to correctly assemble repeated sequences and many so-called whole-genome sequences were in fact incomplete. For example, about 8% of the human genome, made up of highly repetitive sequences, remained unknown

for two decades, until the application of special technologies this year. Similarly, the study of copy number variations due to segmental duplications, which have long been underestimated, is only just beginning. And let us not forget that our exploration of the living world is not only far from complete but is very much biased in favor of the already well-known groups of organisms. We can therefore expect new discoveries, or even surprises, in the study of this part of genomes, overlooked for too long, which demonstrates just how real long-term success differs from the illusion of immediate performance.

Gif-sur-Yvette
Bernard DUJON
Professor Emeritus at *Sorbonne Université*
and the *Institut Pasteur*
Member of the *Institut de France*
(*Académie des Sciences*)
May 2023

Introduction

About Repeated Genomes

Guy-Franck RICHARD

*Instabilités naturelles & synthétiques des génomes,
Institut Pasteur, CNRS UMR3525, Paris, France*

Genome ['dʒi:nəʊm] *nm Biol.*: Set of hereditary characteristics of a living being, of which a small part is composed of genes providing a function to the organism, and the majority is composed of repeated sequences for which it is unknown whether they have a function.

Taking matters a little further, this could be a modern definition of the word “genome”, in the light of the knowledge garnered across three decades from sequencing the DNA content of living beings, in particular eukaryotic organisms, more complex than those of their bacterial and archaeobacterial ancestors. Biologists were already aware back in the 1960s, long before the invention of the first DNA sequencing methods, that the content of genomes was difficult to comprehend. Denaturation-renaturation experiments highlighted that the speed of renaturation of the double-helix was proportional to its concentration. The C_0t parameter was the value at which renaturation of half the genomic DNA was complete, under controlled conditions. Each organism could then be defined by the C_0t value of its genome. In trying to establish the C_0t values of genomes of the simplest organisms - phages or bacteria - or of more complex organisms, such as vertebrates, it transpired that the latter contained three types of sequences presenting very different C_0t values (see [Figure I.1](#)).

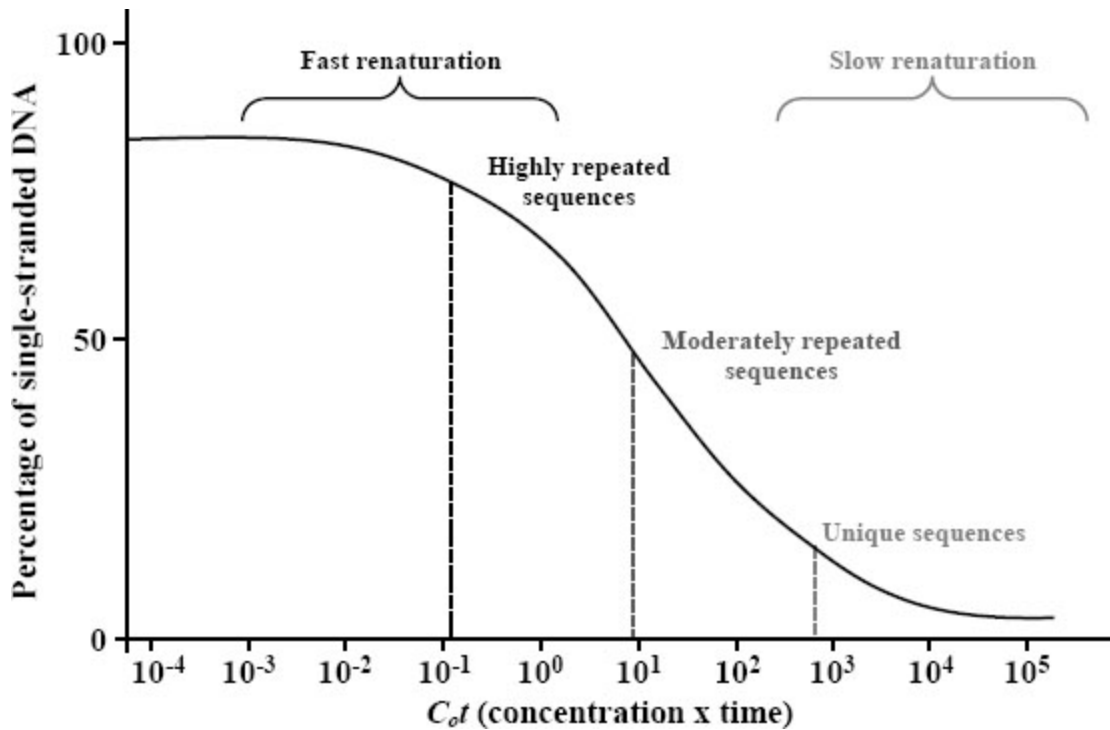


Figure I.1. Example of C_0t curve

It is thus possible to show that the mouse genome, for example, is composed of 70% unique sequences with slow renaturation, 20% moderately repeated sequences present in 1,000 to 100,000 copies per genome and 10% highly repeated sequences representing at least 1 million copies per genome and showing rapid renaturation (Britten and Kohne 1968). This approach, based on the physicochemical properties of DNA, slightly underestimated the quantity of repeated sequences because their renaturation rate depends on the identity between these sequences, divergent sequences (such as long terminal repeats (LTRs)) renaturing more slowly than identical sequences. Nowadays, C_0t curves are still sometimes used to separate the highly repetitive fraction of a genome from its unique fraction in order to sequence specific DNA of either fraction (Peterson et al. 2008).

I.1. The “C-value” paradox

From the moment it was proven that DNA was the support of heredity, and theoretically contained all the genes necessary for the development of a living being, it seemed logical that the most sophisticated organisms had to contain more genes and therefore more DNA in their genome (the “C value”) to encode these genes. This idea was to be questioned in the 1950s with the discovery that the nuclei of certain amphibians and fish contained 20 times more DNA than the nuclei of mammals. Given that the latter presented a greater developmental complexity, this appeared very much paradoxical, and was even used as an argument by the opponents of DNA being the sole support of heredity (Thomas 1971). This “C-value paradox” could finally be explained only decades later, when the first genomes were sequenced. It is now known that the number of genes in an organism has little to do with its size or level of complexity. The baker’s yeast genome contains about 6,000 genes, that of fruit flies about 14,000 and the human genome (or those of its very close cousins, great apes) contains itself with 20,000 genes, with which it manages a very sophisticated level of developmental and behavioral complexity. But what about the paramecium with its 40,000 genes, twice as many as the human genome? Or *Trichomonas vaginalis*, a parasite of the genital tract, with its 60,000 genes? Or indeed wheat and its 124,000 genes, more than six times as many as our genes? Clearly, this so-called complexity could not be measured by the number of genes in an organism. Studies of comparative genomics¹ have shown that this high number of genes in certain organisms does in fact conceal ancestral events of partial or total genome duplication, followed by variable amounts of gene losses (Wolfe and Shields 1997; Jaillon et al. 2004). These events actively participate in the genetic redundancy and their identification as well as their underlying mechanisms will be addressed in [Chapter 1](#).

If the complexity of an organism has nothing to do with the number of genes contained, the same is true of the amount of DNA. The human genome, with just over three times as many genes as brewer's yeast, contains 200 times more DNA. The genome of a rotifer – a small animal measuring just a few millimeters that lives in freshwaters – contains three times more genes than the human genome in 12 times less DNA! (see [Figure I.2](#)).

The genomic sequence of all these organisms showed that some of them had evolved a very compact genome, with high gene density, while others contained a multitude of repeated DNA sequences whose function did not appear obvious at first glance, and that some authors did not hesitate to call them “junk DNA” (Ohno 1972).

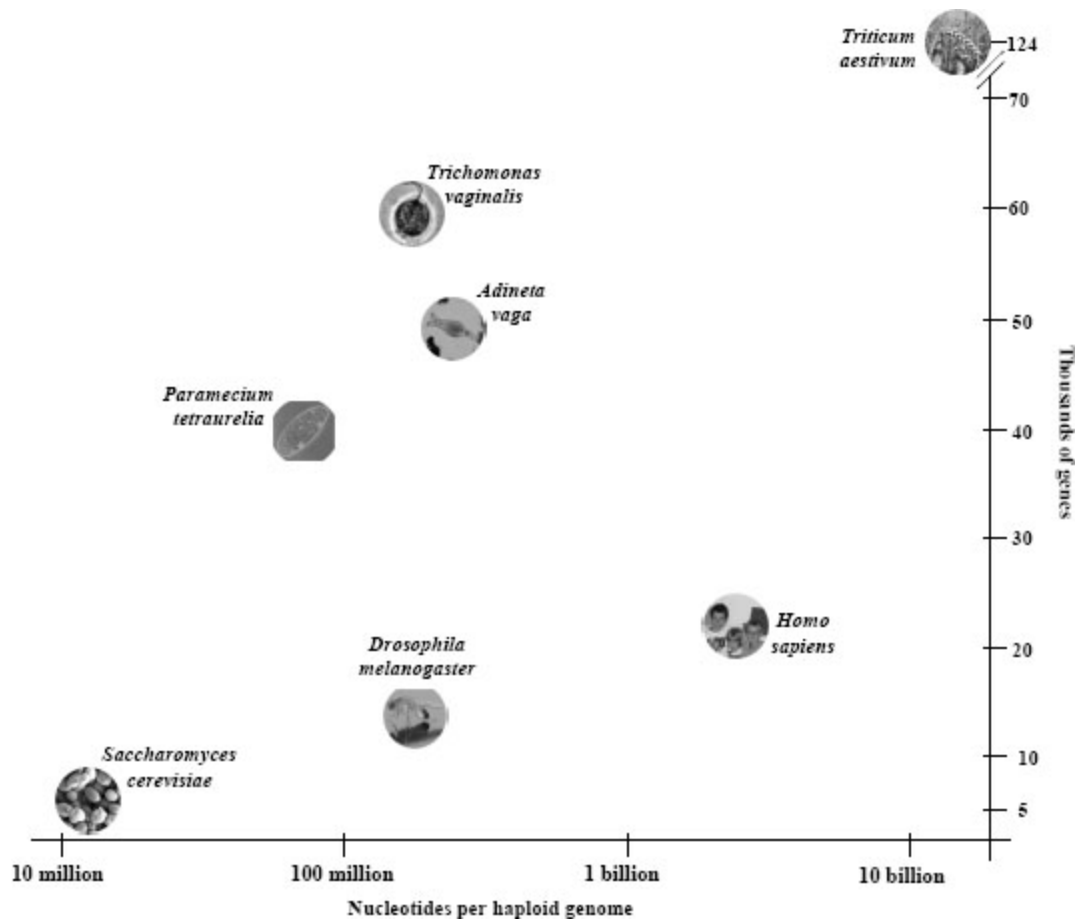


Figure I.2. Comparison of genome sizes and gene numbers

I.2. Recycling junk DNA

About 2% of the human genome is translated into proteins. Even by adding the untranslated genes (rRNA, tRNA, siRNA, snRNA, etc.), the percentage of “useful” DNA barely increases. So, what is the purpose of the 98% of DNA in our genome that has, apparently, no function? One conceivable answer is that it has none. The consortium led by Jeff Boeke, professor of genetics at Johns Hopkins University in Baltimore, set out to create the first synthetic yeast genome, using synthetic oligonucleotides. The brewer’s yeast *Saccharomyces cerevisiae* is a eukaryotic organism whose genome contains 12.5 million nucleotides distributed across 16 chromosomes. The synthetic chromosomes were reconstructed one by one from 70 nucleotide-long sequences assembled in blocks of 750 base pairs, themselves assembled in mega-blocks of 2–4 kb, reintroduced one after the other in a hierarchical manner into the yeast genome in replacement of the natural sequences (Muller and Koszul 2015). When designing synthetic chromosomes, it was decided that all repeated sequences would be removed from the genome. All tRNA-encoding DNAs were grouped on a single circular chromosome, specifically built to carry them. Retrotransposons, microsatellites, minisatellites and other repeated elements inessential to life were removed from the new sequence. These synthetic chromosomes, with their junk DNA removed, are perfectly able to sustain life in yeast cells containing them, without any apparent phenotypic defect, at least under laboratory growth conditions (Dymond et al. 2011; Annaluru et al. 2014). One may conclude from the results of this project that junk DNA is useless. However that would be a mistake.

The human reference genome contains about 443,000 residual elements of past retroviral invasions, covering

8.3% of the total sequence (International Human Genome Sequencing Consortium 2001). These retroviral scars are the remains of successive invasions, occurring over the past hundred million years, of our mammalian ancestors by exogenous elements, which left the trace of their passage in the form of LTR². These retroviral remains are therefore part of our junk DNA. Nevertheless, as we will see, their presence in our genome testifies to their distant but indispensable role in the existence of our lineage. Therian mammals, that is, those possessing a uterus within which the fertilized egg develops, are classified into two groups. Eutherians (or placentals) like humans and mice have a very elaborate placenta connecting the wall of the uterus to the embryo and allowing it to develop in complete safety throughout the entire gestation period. Marsupials (kangaroos and koalas) do not have placentas and the development of their young takes place mainly outside the uterus. Genome sequencing showed that the two human genes specifically expressed in the placenta, *syncytin-1* and *syncytin-2*, were derived from a gene encoding an ancestral viral protein, which infected the primate lineage 25–40 million years ago. Remarkably, the genome of the mouse, another placental mammal, also contains two viral genes having the same function as human genes but deriving from a slightly more recent viral infection than that of the human lineage. Thus, the placenta was invented twice, independently, in two lineages of mammals, by capture of genes of retroviral origin (Dupressoir et al. 2009). Another example is even more striking. Sexual reproduction was invented at the origin of the eukaryotic world. From the first primitive eukaryotic cells, a syngamy³ system was developed that allowed the nuclei of two haploid cells to fuse to give birth to a diploid cell. The protein responsible for the fusion of male and female gametes is the same in plants and animals; it is the product of the *HAP2* gene. This

protein is of viral origin and allows the envelope of a virus to fuse with the plasma membrane of its host's cells (Fédry et al. 2017). Thus, a gene essential to sexual reproduction was captured from a virus by the genome of the very first eukaryotic cells about 1.5 billion years ago.

Other examples of the capture of a piece of transposable element exist, thus creating a new gene, a new function. Junk DNA is therefore regularly recycled during the course of evolution to bring diversity and novelty. As François Jacob (1977) said more than 40 years ago, evolution “tinkers”, it makes new from old, reusing bits of genes, cutting them, splicing them and fusing them with others in order to create novelty. What appears today to the geneticists of the 21st century as junk DNA perhaps served in the past - or will serve in the future - to create diversity. The tremendous success of the eukaryotic world in invading all ecological niches under all climates and latitudes stems in part from the extraordinary flexibility of its genome and its ability to accumulate genetic elements that are seemingly useless but will be recycled in the long run to create novelty and enable the appearance of new living species.

1.3. The different repeat types

There are often several ways to classify genetic elements. Some authors have chosen to distinguish between dispersed repeated elements in contrast to tandem repeats, the latter being repeated at least twice in a row at the same genetic locus, unlike the former, which are repeated at different loci (Richard et al. 2008). But some dispersed repeats are so numerous in the genomes that they appear to be tandemly repeated. This is the case for *Alu* sequences in humans, which are frequently found grouped in introns or intergenic sequences. Repeated sequences of exogenous

origin, that is to say originating from an organism other than the cell in which they are observed, could also be distinguished from repeated sequences of endogenous origin, manufactured by the cell in which they are observed. Transposable elements would belong to the first category, having invaded the genomes of eukaryotic (or prokaryotic) lineages, while the different satellite DNAs would belong to the second, being manufactured by molecular processes specific to the genomes that contain them. But other problems then arise. It is known, for example, that *Alu* elements, inactive retrotransposons that can be mobilized in *trans* by the machinery of other retroelements, are of endogenous origin. They result in fact from the duplication of the non-coding 7SL RNA, which is involved in the synthesis of excreted proteins. This duplication, prior to mammalian radiation, resulted in the fusion of two monomers of 130 nucleotides derived from 7SL RNA, separated by a short adenine-rich region (Ullu and Tschudi 1984). Achieving a coherent classification of the repeated sequences therefore proves a complicated task, particularly in the genomes of evolved plants and animals within which they are plethoric, both in structure and number.

We have therefore tried in the rest of this work to present the repeat elements in relation to their role (proven or assumed) in genomes, rather than according to their structure or their assumed origin (see [Figure I.3](#)).

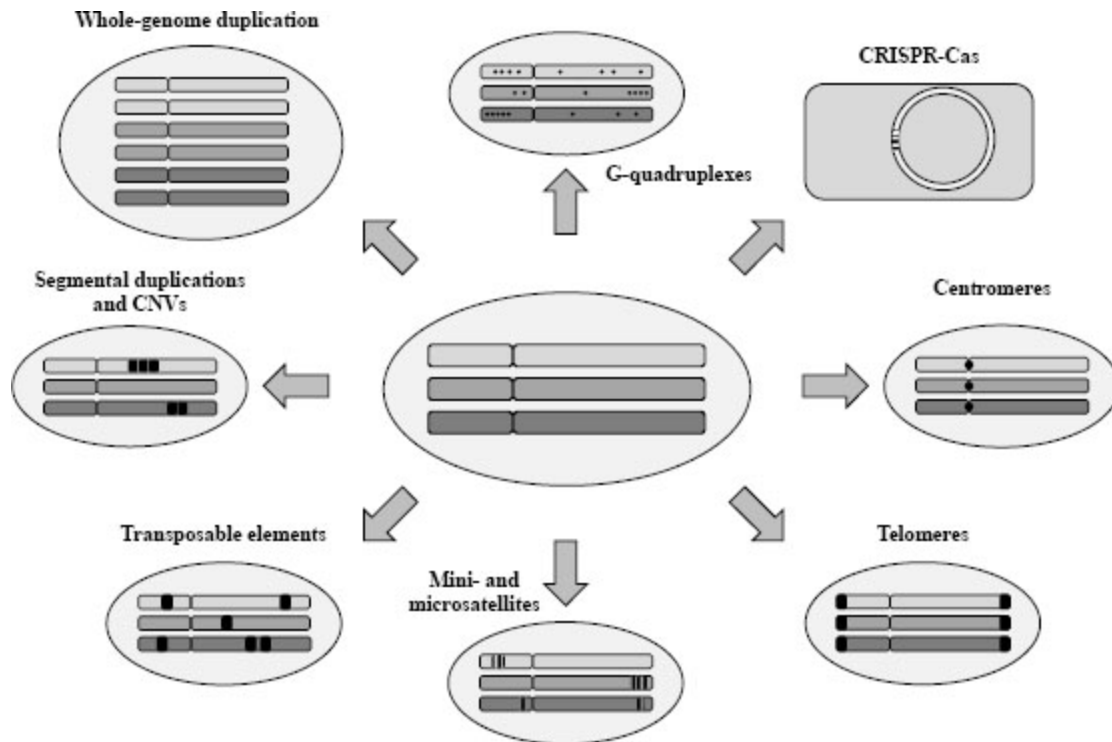


Figure I.3. *The different types of repeated DNA sequences*

After exploring total or partial genome duplications in [Chapter 1](#), the duplications of large DNA segments, sometimes in multiple copies in tandem or dispersed within genomes, will be described in [Chapter 2](#). These contribute significantly to the level of genetic redundancy and gene duplication and their study, although essential to understand the dynamics of complex genomes and the inheritability of certain traits is still in its infancy.

Transposons and retrotransposons will be presented in [Chapter 3](#), and their role in the generation of genetic novelties will be detailed. In most species, centromeres are present at a rate of one per chromosome. These very particular repeated elements are essential for the proper segregation of sister chromatids during cell divisions. They will be studied in [Chapter 4](#) and as we will see, holocentric organisms depart from this rule by exhibiting several tens of centromeres per chromosome. Telomeres are highly repeated sequences found at the ends of chromosomes to

DNA (*see also* [slippage](#))

B-, [245](#), [252](#)

cleavage, [341](#)

junk, [xvii](#)-xx

non-canonical, [239](#), [245](#), [246](#), [249](#), [252](#), [253](#), [255](#)-257

satellite, [273](#)-275, [279](#), [288](#), [289](#)

duplicon, [49](#), [55](#), [95](#)

end-replication problem, [219](#), [220](#)

endogenous retrovirus (ERV), [119](#), [126](#), [130](#), [133](#), [149](#), [150](#)

epigenetics, [124](#), [129](#), [132](#), [138](#), [186](#), [190](#)

evolution, [181](#), [186](#), [187](#), [190](#), [192](#), [193](#)

evolutionary innovations, [23](#)-25, [32](#), [34](#)

forensic research, [278](#)

fragile sites, [299](#), [300](#)

functional innovations, [32](#)

fungi, [184](#), [189](#), [190](#)

G, H

G-quadruplex, [239](#)-250, [252](#)-264

gene

amplification, [49](#), [89](#), [95](#)

conversion, [48](#), [57](#), [60](#), [65](#), [71](#), [72](#), [95](#)

genetic, [181](#), [183](#), [186](#), [189](#), [190](#), [196](#)

duplications, [2](#), [5](#), [18](#)

stability, [258](#)

genome (*see also* [genomic](#)), [xv](#)-[xxii](#)
 evolution, [117](#), [118](#), [140](#), [144](#)
 prokaryote, [319](#), [321](#), [326](#), [331](#), [340](#), [341](#)
 whole-genome duplications, [1](#)-[5](#), [7](#), [9](#), [11](#), [13](#), [15](#)-[18](#), [21](#),
 [22](#), [24](#), [25](#), [27](#), [29](#), [30](#), [32](#)-[35](#)
genomic (*see also* [genome](#)), [259](#), [262](#)
genotyping, [277](#)-[281](#)
holocentromere, [184](#), [192](#), [193](#), [195](#), [197](#), [198](#)
homologs, [2](#), [15](#), [18](#), [19](#), [35](#)
housekeeping genes (HKG), [76](#)

I, L

insects, [193](#), [196](#), [197](#)
interference, [328](#), [334](#)-[343](#), [348](#), [351](#)
lagging strand, [216](#)-[220](#)
leading strand, [217](#)-[220](#)
long terminal repeats (LTRs), [xvi](#), [xix](#)
low-copy repeat (LCR), [58](#), [83](#), [95](#)

M, N

meiosis, [5](#), [7](#)-[10](#), [15](#), [16](#), [19](#)-[21](#)
microsatellite, [48](#), [50](#), [71](#), [72](#), [83](#), [95](#), [273](#)-[287](#), [289](#)-[292](#),
[294](#), [296](#)-[304](#), [307](#)-[311](#)
minisatellite, [48](#), [50](#), [71](#), [83](#), [84](#), [95](#), [273](#)-[280](#), [282](#), [285](#)-[287](#),
[293](#), [294](#), [296](#), [300](#), [304](#), [305](#), [311](#)
molecular domestication, [148](#)

monocentromere, [183](#), [184](#)
mutagenesis, [140](#), [141](#)
nematodes, [182](#), [183](#), [193](#)
next-generation sequencing, [249](#), [250](#)
non-allelic homologous
 recombination (NAHR), [65](#)

P

paralogs, [2](#), [18](#), [22](#), [23](#), [27](#)
paternity test, [278](#)-280
penetrating, [78](#), [95](#)
phage, [321](#), [322](#), [326](#), [331](#)-334, [339](#)-343, [345](#), [346](#), [348](#)-
352
point centromere, [184](#), [186](#), [189](#), [195](#), [198](#)
polyploidies, [4](#), [6](#), [7](#), [35](#)
positive selection, [60](#), [77](#), [80](#), [90](#), [93](#), [95](#)
pseudogene, [74](#), [95](#)
purifying selection, [52](#), [59](#), [62](#), [75](#), [93](#), [95](#)

R, S

redundancy, [1](#), [5](#), [6](#), [9](#), [13](#), [15](#), [16](#), [22](#)-24, [33](#), [34](#)
replicative helicase, [217](#), [218](#)
restriction fragment length polymorphism (RFLP), [277](#), [278](#)
retrotransposons, [120](#)-130, [132](#), [139](#), [142](#)-145, [147](#), [148](#),
[151](#), [154](#)
reverse transcriptase, [209](#), [220](#), [221](#)

RNA (*see also* [CRISPR RNA](#))

small non-coding, [134](#)

structure, [258](#), [259](#)

segmental duplication (SD), [47](#)-50, [52](#), [53](#), [69](#), [95](#)

slippage

during DNA repair, [292](#)

during homologous recombination, [295](#)

during replication, [283](#), [290](#), [297](#), [304](#)

structural variants, [60](#), [95](#), [96](#)

T, V

T-loop, [213](#), [215](#), [216](#), [218](#)

telomerase, [208](#), [209](#), [214](#), [220](#)-224, [226](#)

telomeres, [207](#)-210, [212](#)-227

translocation, [61](#), [63](#), [96](#)

transposition mechanisms, [120](#), [121](#), [123](#), [152](#)

transposons, [xxi](#), [120](#), [121](#), [123](#)-128, [132](#), [141](#), [150](#)-154

trinucleotide repeats, [277](#), [279](#), [281](#), [283](#)-285, [291](#), [292](#), [296](#), [297](#), [299](#)-301, [304](#)

virus, [xix](#)