



SCIENCES

BIOLOGY

Genetics, Epigenetics

Function and Evolution of Repeated DNA Sequences

**Coordinated by
Guy-Franck Richard**

ISTE

WILEY

Function and Evolution of Repeated DNA Sequences

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

Contents

Foreword	xiii
Bernard DUJON	
Introduction	xv
Guy-Franck RICHARD	
Chapter 1. Whole-Genome Duplications, a Source of Redundancy at the Entire-Genome Scale.	1
Elise PAREY and Camille BERTHELOT	
1.1. Prevalence of polyploids in the tree of life	2
1.1.1. Whole duplications in eukaryotes	2
1.1.2. Polyploidies in prokaryotic organisms.	6
1.1.3. Polyploid cells in normal and pathological physiology.	7
1.2. Mechanisms for the appearance of whole-genome duplications.	7
1.2.1. Non-separation of chromosomes after replication.	7
1.2.2. Autopolyploidization, a perfect genome redundancy	9
1.2.3. Allopolyploidization, an overlapping of genomes of similar species	9
1.3. Cellular consequences of whole-genome duplications	11
1.3.1. Disruption of cell and nucleus organization	11
1.3.2. Modifications in the expression of genes and transposons	13
1.3.3. Unstable meiosis	15
1.4. Rediploidization: evolutionary reduction in genetic redundancy	16
1.4.1. Resolution of meiosis by karyotype rearrangement	16

1.4.2. Evolutionary divergence of duplicated sequences	18
1.4.3. Bias and dominance during rediploidization	20
1.4.4. Incomplete and lineage-specific rediploidizations.	21
1.5. Functions and evolution of duplicated genes	22
1.5.1. Redundancy and subfunctionalization	23
1.5.2. Neofunctionalization and evolutionary innovations.	24
1.5.3. Gene repertoire bias	26
1.5.4. Regulatory blocks and splitting of regulatory regions	29
1.6. Whole-genome duplications and evolutionary diversification	32
1.6.1. Association with geological crises	32
1.6.2. Evolutionary speciations and radiations.	33
1.7. Perspectives and conclusions	34
1.8. References	35
Chapter 2. Segmental Duplications and CNVs: Adaptive Potential of Structural Polymorphism	47
Patricia BALARESQUE and Franklin DELEHELLE	
2.1. The multiple facets of genetic polymorphism.	48
2.2. From Segmental Duplications to Copy Number Variants: terminology	49
2.3. SDs: a general overview	49
2.3.1. Background	49
2.3.2. SDs: more than a category of sequences, superstructures	50
2.3.3. SD and CNV: study biases related to the attractiveness of subjects as well as to the technological developments of the moment	51
2.3.4. SD: characteristics in human and non-human primates.	52
2.4. Methodologies for detecting structural variation in genomes	53
2.4.1. In vitro methods.	54
2.4.2. Methods on reads	54
2.4.3. Post-assembly methods	54
2.5. The molecular mechanisms at the origin of structural variation	56
2.5.1. Homologous recombination mechanisms	56
2.5.2. Non-homologous recombination mechanisms	57
2.6. Regions rich in SDs/LCRs favor the creation of CNVs: insertions/duplications, deletions and inversions	58
2.6.1. Insertions/duplications and deletions	58
2.6.2. Inversions	60
2.7. From SDs to CNVs in humans and primates	61
2.7.1. General overview	61
2.7.2. Delineating regions of interest	61

2.7.3. Heterogeneity in the distribution of intra- and interchromosomal SDs.	62
2.7.4. Intrachromosomal and interchromosomal SDs: what do they teach us about the evolutionary history and origin of SDs?	62
2.7.5. Intra- and interchromosomal SDs: the specific case of sex chromosomes.	66
2.7.6. SDs: an association with specific sequences?	66
2.8. SDs in little-studied species: general genomic profiles	66
2.8.1. Twelve genomes under study	68
2.8.2. Distribution and characteristics of SDs in genomes.	70
2.9. SD content: impact of a duplicated environment on sequences that make up the SDs.	70
2.9.1. SDs and non-coding sequences: the case of microsatellites	71
2.9.2. SDs and coding genes: the fate of genes in SDs	72
2.10. SDs and epigenetic modifications	75
2.11. The adaptive potential of SDs: between the benefit of innovation and the cost of pathology	78
2.11.1. The organism's defense: immune system	79
2.11.2. Nutrient/food assimilation	80
2.11.3. Sensory perception of the environment	80
2.11.4. Neurological processes	82
2.11.5. Reproduction and the X and Y chromosomes: true SD concentrates	83
2.12. SDs and associated CNVs: their roles in species adaptation to changes in environments.	86
2.12.1. SDs: a link between genomic architecture, adaptive potential and environmental changes?	86
2.12.2. Adaptation to global environmental stress.	86
2.12.3. Adaptation to nutrient-poor surroundings	88
2.12.4. Adaptation to low and high temperatures	88
2.12.5. Heavy-metal adaptation	89
2.12.6. Antibiotics and drugs	90
2.12.7. Pesticide resistance	90
2.12.8. Domestication and post-domestication of plant and animal species.	91
2.12.9. Competition and evolutionary success: invasive species and hybridization	93
2.13. Conclusion	94
2.14. Glossary of terms.	95
2.15. References.	96

Chapter 3. Transposable Elements: Parasites that Shape Genome Evolution 117

Amandine BONNET, Karine CASIER, Clément CARRÉ,
Laure TEYSSET and Pascale LESAGE

3.1. Transposable elements in eukaryotic genomes	117
3.1.1. TEs: essential components of eukaryotic genomes	118
3.1.2. Acquisition of new TEs by horizontal transfer	119
3.2. Classification of TEs and transposition mechanisms	120
3.2.1. Class I retrotransposons	120
3.2.2. Class II DNA transposons	123
3.3. TE self-regulation	123
3.3.1. Spatio-temporal regulation of TE expression.	124
3.3.2. Self-regulation of transposition efficiency	125
3.3.3. Selective integration to better protect the genome.	125
3.4. TE restriction by the host.	129
3.4.1. Transcriptional repression of genomic copies	129
3.4.2. TE transcripts: choice targets for multiple restrictions	132
3.4.3. The Swiss knives of TE restriction: piRNAs	134
3.4.4. Reverse transcription of retroelements: a key step to inhibit.	139
3.5. The impact of transposition events on genomes	140
3.5.1. The structural and functional consequences of TE activity on the genome.	140
3.5.2. Pathologies associated with TE activity.	144
3.5.3. The impact of TEs on the evolution of the host	148
3.6. Conclusion	155
3.7. References	156

Chapter 4. Insights Into the Evolutionary Diversity of Centromeres. 181

Nuria CORTES-SILVA, Aruni P. SENARATNE and Ines A. DRINNENBERG

4.1. The centromere.	181
4.1.1. Definition and historical background	181
4.1.2. Two main types of centromeric architectures	183
4.2. Monocentromeres	184
4.2.1. The diversity of monocentric architectures across fungi	184
4.2.2. Animal and plant models contain long repetitive regional centromeres.	190
4.3. Holocentromeres.	192
4.3.1. Nematodes	193

4.3.2. Plants	195
4.3.3. Insects.	196
4.4. Open questions.	198
4.5. Acknowledgments.	198
4.6. References	198
Chapter 5. Evolution and Functions of Telomeres	207
Arturo LONDOÑO-VALLEJO	
5.1. Primary structure of telomeres.	207
5.1.1. Origin and evolution of telomeres	210
5.1.2. Nucleoprotein structure of telomeres	212
5.2. A telomere specific higher order structure: the T-loop	215
5.2.1. Telomere replication, a fundamental mechanism for telomere maintenance	215
5.3. Telomere lengthening mechanisms	220
5.4. Telomere length homeostasis	222
5.5. Telomeres and genome organization and function	225
5.6. Cell senescence, aging and disease	226
5.7. Conclusion	227
5.8. Acknowledgments.	227
5.9. References	227
Chapter 6. G-quadruplexes: Structure, Detection and Functions	239
Emilia Puig LOMBARDI	
6.1. From guanine-guanine base-pairing to a secondary structure	239
6.1.1. G-quartets	239
6.1.2. Folding into a G-quadruplex structure.	241
6.2. The G4 structure: variations on a theme	243
6.2.1. RNA G-quadruplexes (rG4).	245
6.2.2. Exceptions to the rule(s): non-canonical G-quadruplexes	245
6.3. Finding G-quadruplexes in a genome	246
6.3.1. Experimental methods for G-quadruplex detection	247
6.3.2. Computational methods	250
6.4. Biological roles of G-quadruplexes.	257
6.4.1. First role attributed to quadruplexes: their formation in telomeres.	257
6.4.2. Predictions based on bioinformatic analyses	259
6.5. Perspective: G-quadruplexes as anticancer therapeutic targets.	261
6.6. References	264

Chapter 7. Satellite DNA, Microsatellites and Minisatellites 273

Wilhelm VAYSSE–ZINKHÖFER and Guy-Franck RICHARD

7.1. Satellite DNAs, origin and definition.	273
7.1.1. Minisatellites.	274
7.1.2. Microsatellites.	274
7.2. From semantics to biology.	275
7.2.1. Distribution of satellite DNAs in genomes	275
7.2.2. Polymorphic genetic markers	277
7.2.3. Trinucleotide repeat expansions	281
7.2.4. Microsatellites regulate gene expression	283
7.2.5. Minisatellites are important in cell adhesion	285
7.2.6. Function of megasatellites.	287
7.2.7. Centromeric satellite DNA, complexity of structure–function studies.	288
7.3. The evolutionary mechanisms of tandem repeats	289
7.3.1. Historical model of slippage during replication	290
7.3.2. Slippage during DNA repair	292
7.3.3. Repeat expansions and contractions during homologous recombination	292
7.4. Microsatellites in human diseases.	297
7.4.1. Triplet repeat expansion disorders	297
7.4.2. Colorectal cancers and the mismatch repair system.	298
7.4.3. Fragile sites	299
7.5. De novo formation and evolution of tandem repeats.	300
7.5.1. Birth and death of microsatellites	300
7.5.2. Formation of minisatellites	304
7.6. Perspectives	307
7.6.1. Inadequacy of software tools	307
7.6.2. The importance of definitions in biology	310
7.7. Acknowledgments.	311
7.8. References	311

Chapter 8. CRISPR-Cas: An Adaptive Immune System 319

Marie TOUCHON

8.1. A brief history of the discovery of CRISPR-Cas systems.	319
8.2. General characteristics of CRISPR-Cas systems	323
8.2.1. Diversity of repeats.	324
8.2.2. Diversity and origin of spacers	325
8.2.3. Diversity and evolutionary classification of cas genes	327

8.2.4. Origin of CRISPR-Cas systems	329
8.3. Evolution of CRISPR-Cas systems	330
8.3.1. Scattered distribution of CRISPR-Cas systems	330
8.3.2. Massive transfer of CRISPR-Cas systems	331
8.3.3. Commonly lost systems	332
8.3.4. Evolutionary dynamics of CRISPR arrays	333
8.4. An adaptive immune system.	334
8.4.1. A three-stage immune response.	334
8.4.2. Diversity of CRISPR-Cas molecular mechanisms.	337
8.4.3. Self- and none self-discrimination: avoiding self-targeting by CRISPR	340
8.5. Phage escape mechanisms	341
8.5.1. Genomic modifications	341
8.5.2. Anti-CRISPR proteins	343
8.6. Biological cost of CRISPR-Cas systems.	344
8.6.1. Cost of expression	344
8.6.2. Cost of autoimmunity	345
8.6.3. The genetic background of the host	346
8.6.4. Limiting horizontal gene transfer.	347
8.6.5. Naïve and primed adaptation	348
8.7. Importance in nature: impact of ecological factors.	349
8.7.1. Phage diversity – mutation rate.	349
8.7.2. Phage diversity – population size.	350
8.7.3. Infectious risk – alternative strategies	350
8.8. Conclusions and perspectives	351
8.9. References	353
List of Authors	361
Index	363

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).

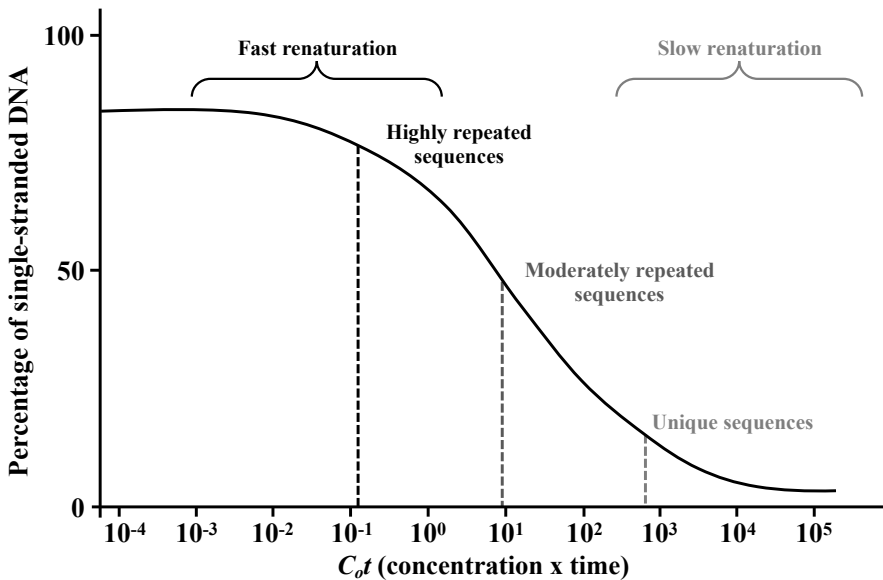


Figure I.1. Example of C_0t curve

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).

¹ Comparative genomics is the field of genomics that focuses on the comparison of entire genomes with each other and not just of genes. Analysis tools have been specifically developed to compare the organization, structure, synteny (gene order along a chromosome) of genomes, considered as objects to be studied as a whole.

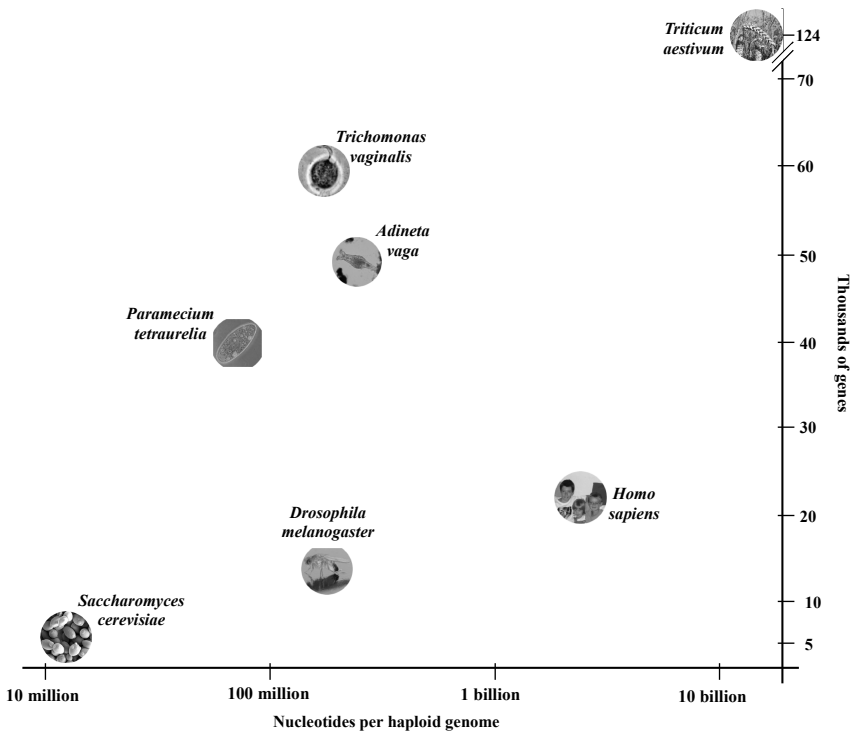


Figure I.2. Comparison of genome sizes and gene numbers

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

² LTR (long terminal repeat): repeated sequences typically found at the insertion sites of retroviruses.

³ Syngamy: nuclear fusion of two cells of opposite mating type.

(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).

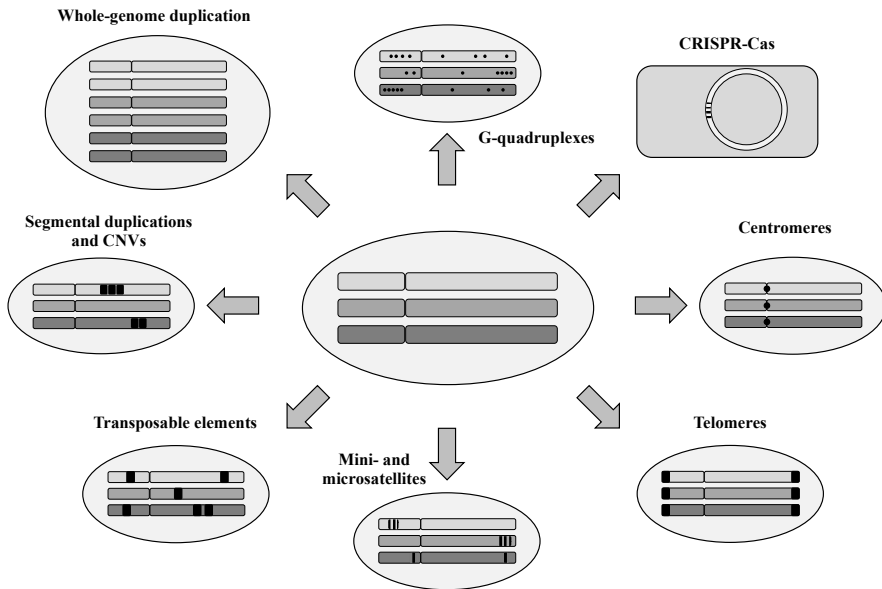


Figure I.3. *The different types of repeated DNA sequences*

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 prevent loss of genetic

information. Their sequence and structure vary greatly from organism to organism, with some species having developed highly original telomeres that are made up of tandemly repeated elements. These concepts will be studied in Chapter 5. G-quadruplexes, these secondary DNA structures caused by the regular repeat of GC base pairs, are present in all eukaryotic genomes. Their distribution and their role in DNA transcription and replication will be discussed in Chapter 6. The different types of satellite DNA found in large numbers in eukaryotes, and whose precise function is not always clear, will be described in Chapter 7. As we will see, although prokaryotic genomes contain only few of them, some bacteria use them as camouflage to escape their host's immune system. Remaining in the world of prokaryotes, we will end in Chapter 8 with the fascinating study of another bacterial defense mechanism, directed against these other enemies that are plasmids and bacteriophages: the CRISPR-Cas system. The acquisition of small, tandemly repeated pieces of DNA from invaders foreign to the cell provides eubacteria and archaea with a robust line of defense. And it offers 21st-century geneticists myriad tools to manipulate their preferred genomes at their own convenience.

I.4. References

- Annaluru, N., Muller, H., Mitchell, L.A., Ramalingam, S., Stracquandano, G., Richardson, S.M., Dymond, J.S., Kuang, Z., Scheifele, L.Z., Cooper, E.M. et al. (2014). Total synthesis of a functional designer eukaryotic chromosome. *Science*, 344(6179), 55–58.
- Britten, R.J. and Kohne, D.E. (1968). Repeated sequences in DNA. *Science*, 161, 529–540.
- Dupressoir, A., Vernochet, C., Bawa, O., Harper, F., Pierron, G., Opolon, P., Heidmann, T. (2009). Syncytin-A knockout mice demonstrate the critical role in placentation of a fusogenic, endogenous retrovirus-derived, envelope gene. *Proceedings of the National Academy of Sciences*, 106(29), 12127–12132.
- Dymond, J.S., Richardson, S.M., Coombes, C.E., Babatz, T., Muller, H., Annaluru, N., Blake, W.J., Schwerzmann, J.W., Dai, J., Lindstrom, D.L. et al. (2011). Synthetic chromosome arms function in yeast and generate phenotypic diversity by design. *Nature*, 477(7365), 471–476.
- Fédry, J., Liu, Y., Péhau-Arnaudet, G., Pei, J., Li, W., Tortorici, M.A., Traincard, F., Meola, A., Bricogne, G., Grishin, N.V. et al. (2017). The ancient gamete fusogen HAP2 is a eukaryotic class II fusion protein. *Cell*, 168(5), 904–915.e10.
- International Human Genome Sequencing Consortium (2001). Initial sequencing and analysis of the human genome. *Nature*, 409, 860–921.
- Jacob, F. (1977). Evolution and tinkering. *Science*, 196(4295), 1161–1166.

- Jaillon, O., Aury, J.-M., Brunet, F., Petit, J.-L., Stange-Thomann, N., Mauceli, E., Bouneau, L., Fischer, C., Ozouf-Costaz, C., Bernot, A. et al. (2004). Genome duplication in the teleost fish *Tetraodon nigroviridis* reveals the early vertebrate proto-karyotype. *Nature*, 431(7011), 946–957.
- Muller, H. and Koszul, R. (2015). Conception et synthèse de néochromosomes. *Médecine thérapeutique/Médecine de la reproduction, gynécologie et endocrinologie*, 17(4), 228–236.
- Ohno, S. (1972). So much “junk” DNA in our genome. *Evolution of Genetic Systems*, 23, 366–370.
- Peterson, D.G., Schulze, S.R., Sciara, E.B., Lee, S.A., Bowers, J.E., Nagel, A., Jiang, N., Tibbitts, D.C., Wessler, S.R., Paterson, A.H. (2008). Integration of Cot analysis, DNA cloning, and high-throughput sequencing facilitates genome characterization and gene discovery. *Genome Research*, 12, 795–807.
- Richard, G.-F., Kerrest, A., Dujon, B. (2008). Comparative genomics and molecular dynamics of DNA repeats in eukaryotes. *Microbiol. Mol. Biol. Rev.*, 72(4), 686–727.
- Thomas Jr., C.A. (1971). The genetic organization of chromosomes. *Annu. Rev. Genet.*, 5, 237–256.
- Ullu, E. and Tschudi, C. (1984). Alu sequences are processed 7SL RNA genes. *Nature*, 312, 171–172.
- Wolfe, K.H. and Shields, D.C. (1997). Molecular evidence for an ancient duplication of the entire yeast genome. *Nature*, 387, 708–713.

Franklin DELEHELLE
Institut de Biologie de l'École Normale
Supérieure
CNRS UMR8197
Université Paris Sciences & Lettres
France

Bernard DUJON
Sorbonne Université
and
Institut Pasteur
Institut de France
Académie des Sciences
Paris
France

Ines A. DRINNENBERG
Evolution des centromères et ségrégation
des chromosomes
Institut Curie
CNRS UMR3664
Paris
France

Pascale LESAGE
Génomés, biologie cellulaire et
thérapeutique
Université Paris Cité
CNRS UMR7212
Inserm U944
France

Emilia Puig LOMBARDI
SFR Necker
CNRS UAR3633
Inserm US24
Paris
France

Arturo LONDOÑO-VALLEJO
Téломères et cancer
Institut Curie
CNRS UMR3347
Paris
France

Elise PAREY
Institut de biologie de l'École Normale
Supérieure
CNRS UMR8197
INSERM U1024
Paris
France

Guy-Franck RICHARD
Instabilités naturelles & synthétiques
des génomes
Institut Pasteur
CNRS UMR3525
Paris
France

Aruni P. SENARATNE
Evolution des centromères et ségrégation
des chromosomes
Institut Curie
CNRS UMR3664
Paris
France

Laure TEYSSET
Epigénétique transgénérationnelle et
biologie des petits ARN
Institut de Biologie Paris Seine
Sorbonne Université
CNRS UMR7622
France

Marie TOUCHON
Génomique évolutive des microbes
Institut Pasteur
CNRS UMR3525
Paris
France

Wilhelm VAYSSE-ZINKHÖFER
Instabilités naturelles & synthétiques
des génomes
Institut Pasteur
CNRS UMR3525
Paris
France

Index

A

adaptation, 327–330, 334–336, 338, 340, 342, 345–352
allopolyploids, 3, 9–11, 13–15, 20, 21, 23
Alu element, xx
autopolyploids, 3, 9, 11, 14, 15

C

capping function, 219
cas9, 327, 328, 338, 339, 343, 353
centromere, 210, 225
colorectal cancers, 298
computational biology, 253
copy number variant (CNV), 49
 C_0t curve, xvi
CRISPR array, 321–323, 325, 327, 333–336, 338, 340, 341, 345, 347, 351
CRISPR RNA, 335

D, E, F

defense system, 319, 327, 340, 341, 348, 352

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