



Huaichun Wang

The Effects of Nucleotide Bias on Genome Evolution

The causes and effects of wide variations in
G+C content of the genomes

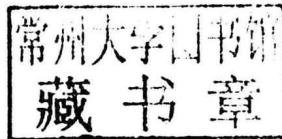
VDM

Verlag
Dr. Müller

Huaichun Wang

The Effects of Nucleotide Bias on Genome Evolution

The causes and effects of wide variations in
G+C content of the genomes



VDM Verlag Dr. Müller

Impressum/Imprint (nur für Deutschland/ only for Germany)

Bibliografische Information der Deutschen Nationalbibliothek: Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über <http://dnb.d-nb.de> abrufbar.

Alle in diesem Buch genannten Marken und Produktnamen unterliegen warenzeichen-, marken- oder patentrechtlichem Schutz bzw. sind Warenzeichen oder eingetragene Warenzeichen der jeweiligen Inhaber. Die Wiedergabe von Marken, Produktnamen, Gebrauchsnamen, Handelsnamen, Warenbezeichnungen u.s.w. in diesem Werk berechtigt auch ohne besondere Kennzeichnung nicht zu der Annahme, dass solche Namen im Sinne der Warenzeichen- und Markenschutzgesetzgebung als frei zu betrachten wären und daher von jedermann benutzt werden dürften.

Coverbild: www.purestockx.com

Verlag: VDM Verlag Dr. Müller Aktiengesellschaft & Co. KG
Dudweiler Landstr. 99, 66123 Saarbrücken, Deutschland
Telefon +49 681 9100-698, Telefax +49 681 9100-988, Email: info@vdm-verlag.de
Zugl.: Ottawa, University of Ottawa, PhD dissertation, 2005

Herstellung in Deutschland:
Schaltungsdienst Lange o.H.G., Berlin
Books on Demand GmbH, Norderstedt
Reha GmbH, Saarbrücken
Amazon Distribution GmbH, Leipzig
ISBN: 978-3-639-19164-6

Imprint (only for USA, GB)

Bibliographic information published by the Deutsche Nationalbibliothek: The Deutsche Nationalbibliothek lists this publication in the Deutsche Nationalbibliografie; detailed bibliographic data are available in the Internet at <http://dnb.d-nb.de>.

Any brand names and product names mentioned in this book are subject to trademark, brand or patent protection and are trademarks or registered trademarks of their respective holders. The use of brand names, product names, common names, trade names, product descriptions etc. even without a particular marking in this works is in no way to be construed to mean that such names may be regarded as unrestricted in respect of trademark and brand protection legislation and could thus be used by anyone.

Cover image: www.purestockx.com

Publisher:
VDM Verlag Dr. Müller Aktiengesellschaft & Co. KG
Dudweiler Landstr. 99, 66123 Saarbrücken, Germany
Phone +49 681 9100-698, Fax +49 681 9100-988, Email: info@vdm-publishing.com

Copyright © 2010 by the author and VDM Verlag Dr. Müller Aktiengesellschaft & Co. KG and licensors
All rights reserved. Saarbrücken 2010

Printed in the U.S.A.
Printed in the U.K. by (see last page)
ISBN: 978-3-639-19164-6

Huaichun Wang

The Effects of Nucleotide Bias on Genome Evolution

To the memory of my father

Acknowledgements

I would like to thank my supervisor, Professor Donal Hickey for assigning me this interesting thesis topic and directing every detail of the whole thesis work, and lots of financial support for attending conferences and summer schools. I also thank the members of my advisory committee, Drs. George Carmody, Guy Drouin and Marcel Turcotte for their advice. Dr. Xuhua Xia also gave very insightful advice that is greatly appreciated. Ada Chyurlia always gave me help in need and Dr. Greg Singer helped a lot by maintaining the Sun server (even after he had left the lab), which was essential for my programming. Financial supports from The Natural Sciences and Engineering Research Council (NSERC) of Canada (to Dr. Hickey), Ontario Graduate Scholarship and University of Ottawa graduate scholarships are appreciated. Finally, as a student at a senior age, family supports are all important. I would like to thank my parents for bringing me up and my brothers and sister for taking care of them so I can study uninterrupted for the past three and a half years. My wife Shuli and my son Sean gave me love and joy so that my study was less stressful.

Abstract

The genomic G+C content of prokaryotes varies from approximately 23% to 77% among genomes. In contrast, among vertebrates, the variation is greatest within the same genome rather than between genomes. There has been a long-standing controversy concerning the causes of these inter- and intra-specific variations. Is it caused by natural selection, favored by the selectionists or, conversely, is it selectively neutral (the neutralist view)? In this study, we investigated the source of nucleotide compositional variation (nucleotide bias) and the consequences of the bias on protein sequence and genome evolution. Thermal adaptation is a primary example to study the effect of natural selection and has been thoroughly studied in this project. We found that both GC content and length of ribosomal RNA genes show positive correlations with optimal growth temperature in prokaryotes and these correlations are not due to phylogenetic history. The correlations are concentrated almost entirely within the stem regions of the rRNA. The rRNA loops, however, show very constant base composition regardless of temperature optima or genomic GC content. The loops were found to have very high amount of adenosine nucleotides throughout prokaryotes and eukaryotes. These results clearly demonstrated that environmental temperature is a selective force that drives rRNA gene evolution and different segments of the same gene (i.e., the stems and loops of the rRNA gene) experience differential selection, although the mutation spectrum presumably should be similar between the loops and stems.

For protein coding genes, mutation and natural selection play a different role compared to the rRNA genes. The neutralist predicts mutational bias would cause protein sequence evolution, while the selectionist would predict that the protein sequence is not related to genomic GC content. To investigate these two postulations and analyze the consequences of nucleotide bias in eukaryotic genomes, we studied homologous genes and their encoded proteins in two flowering plants, *Oryza sativa* (rice) and *Arabidopsis thaliana*. While there is a relatively homogenous GC content in the *Arabidopsis* genes (26% to 69%), the GC content of the rice genes is very heterogeneous (27% to 90%). High GC rice genes encode proteins having a high frequency of GC-rich codons encoded amino acids, i.e., glycine, alanine, arginine and proline. Low GC rice genes and

Arabidopsis genes encode proteins having a high frequency of AT-rich codons encoded amino acids, *i.e.*, phenylalanine, tyrosine, methionine, isoleucine, asparagines and lysine. Furthermore, the effects of nucleotide bias on synonymous codon usage in the rice and *Arabidopsis* genomes were studied. We have shown that synonymous codon usage in the rice genome is primarily dictated by the GC content of the genes, rather than by translational selection. This study in multicellular higher plants, together with previous work on prokaryote and yeast, provide persuasive evidence that mutational nucleotide bias is a cause, rather than a consequence, of protein evolution and this affects codon usage and protein composition in a predictable way.

Résumé

Le contenu génomique en nucléotides G+C des prokaryotes varie approximativement de 23% à 77% entre génomes. Chez les vertébrés, contrairement aux prokaryotes, cette variation est plus élevée au sein d'un génome plutôt qu'entre génomes différents. Il existe une controverse de longue durée concernant les causes de ces variations inter- et intra-spécifique. Est-ce que ces variations sont le résultat de la sélection naturelle, cette explication est favorisée par les sélectionnistes ou d'une sélection neutre (point vue des neutralistes)? Dans cette étude, nous examinons la source des variations de la composition en nucléotides (biais nucléotidique) et les conséquences de ce biais sur les séquences de protéines et sur l'évolution du génome. L'adaptation thermique convient bien à l'étude de l'effet de la sélection naturelle et nous l'avons étudié de façon rigoureuse dans ce projet. Nous avons trouvé que le contenu en nucléotides GC et la longueur des gènes d'ARN ribosomal (ARNr) montrent une corrélation positive avec les températures optimales de croissance chez les prokaryotes et que ces corrélations ne sont pas le résultat de l'histoire phylogénétique des espèces. Les corrélations sont concentrées principalement au niveau des tiges des ARNr. Par contre, les bras d'ARNr sont très constants dans leur composition en nucléotides et ne sont pas affectés par les températures optimales ou le contenu en nucléotides GC. Il semblerait que les bras abondent de nucléotides d'adénosines et ce autant chez les prokaryotes que chez les eucaryotes. Ces résultats démontrent clairement que la température environnementale exerce une force sélective qui conduit à l'évolution des gènes ARNr et que différents segments du même gène (i.e., les tiges et les bras du gène d'ARNr) sont affecté par différente sélection même si le spectre des mutations est présumé être similaire dans les tiges et les bras.

Pour les gènes codant pour des protéines, les mutations et la sélection naturelle jouent un rôle différent comparé aux gènes ARNr. Les neutralistes prédissent qu'un biais mutationnelle chez les séquences protéiques n'est pas associé au contenu en GC génomique. L'investigation de ces deux postulats et analyser les conséquences du biais nucléotidique, nous avons étudié des gènes homologues et qui codent pour des protéines de deux plantes à fleurs: *Oryza sativa* (riz) et *Arabidopsis thaliana*. Pendant que le contenu en GC est relativement homogène dans les gènes d'*Arabidopsis* (26% à 69%), le

contenu en GC chez le riz est très hétérogène (27% à 90%). Les gènes de riz à haute teneur en GC codent pour des protéines dont la fréquence d'acides aminés ayant des codons riches en nucléotides GC est élevée, i.e. glycine, alanine, arginine et proline. Chez le riz et *Arabidopsis*, les gènes à plus faible teneur en GC codent pour des protéines composées d'acides aminés ayant des codons riches en base AT, i.e. phenylalanine, tyrosine, methionine, isoleucine, asparagine et lysine. De plus, les effets du biais nucléotidique sur l'utilisation des codons synonymes du génome du riz est principalement contrôlé par le contenu en nucléotides GC des gènes, plutôt que par la sélection traductionnelle. Cette étude des plantes multicellulaires d'ordre supérieur, de même que des recherches passées sur les prokaryotes et levures, démontrent de façon claire et précise que le biais nucléotidique est une cause plutôt qu'une conséquence de l'évolution des protéines et ceci affecte l'utilisation des codons et la composition protéique que nous pouvons prédire.

RC: (codon) redundancy class

rRNA: ribosomal RNA

RSCU: Relative Synonymous Codon Usage

ssu rRNA: small subunit rRNA

Thermo: thermophiles (thermophilic species)

tRNA: transfer RNA

TT: thymine dimer

UV: ultraviolet

VSP: (DNA repair) very short patch

IUPAC code table

Amino acid codes			Nucleic acid codes	
1-letter	3-letter	description	code	description
A	Ala	Alanine	A	Adenine
R	Arg	Arginine	C	Cytosine
N	Asn	Asparagine	G	Guanine
D	Asp	Aspartic acid	T	Thymine
C	Cys	Cysteine	U	Uracil
Q	Gln	Glutamine	R	Purine (A or G)
E	Glu	Glutamic acid	Y	Pyrimidine (C, T, or U)
G	Gly	Glycine	M	C or A
H	His	Histidine	K	T, U, or G
I	Ile	Isoleucine	W	T, U, or A
L	Leu	Leucine	S	C or G
K	Lys	Lysine	B	C, T, U, or G (not A)
M	Met	Methionine	D	A, T, U, or G (not C)
F	Phe	Phenylalanine	H	A, T, U, or C (not G)
P	Pro	Proline	V	A, C, or G (not T, not U)
S	Ser	Serine	N	Any base (A, C, G, T, or U)
T	Thr	Threonine		
W	Trp	Tryptophan		
Y	Tyr	Tyrosine		
V	Val	Valine		
B	Asx	Aspartic acid or Asparagine		
Z	Glx	Glutamine or Glutamic acid		

Table of Contents

Acknowledgements.....	iii
Abstract.....	iv
Résumé.....	vi
List of Abbreviations.....	viii
IUPAC Code Table.....	x
List of Tables.....	xv
List of Figures.....	xvii
Chapter 1 General introduction.....	1
1.1 Variations (bias) in \bar{GC} content.....	1
1.2 Effects of GC variations.....	2
1.3 Biased DNA mutation and natural selection shape \bar{GC} content.....	5
1.3.1 Selectionist interpretations of GC content change.....	5
1.3.2 Neutralist interpretations of genomic differences in GC content.....	7
1.3.3 DNA mutation.....	9
1.3.4 Bias in mutation.....	11
1.3.5 DNA mutation repair and repair bias.....	12
1.3.6 Interaction of neutral evolution and selection.....	13
1.4 Research proposal.....	14
1.5 Comparative methods.....	15
1.6 Organization of the thesis.....	18
Chapter 2 Evidence for strong selective constraint acting on the nucleotide composition of 16S ribosomal RNA genes.....	21
2.1 Abstract.....	21
2.2 Introduction.....	22
2.3 Materials and Methods.....	23
2.4 Results.....	24

2.4.1	Average nucleotide composition in mesophiles and thermophiles.....	24
2.4.2	16S rRNA stems and loops are affected very differently by growth temperature.....	28
2.4.3	The relationship between the nucleotide content of the 16S rRNA and the nucleotide content of the whole genome.....	31
2.5	Discussion.....	34
Chapter 3	Thermal adaptation of ribosomal RNA genes.....	39
3.1	Abstract.....	39
3.2	Introduction.....	40
3.3	Methods.....	42
3.3.1	Sequence data.....	42
3.3.2	Growth temperature.....	42
3.3.3	Statistical analyses.....	44
3.4	Results.....	46
3.4.1	Nucleotide composition and sequence length of 16S rRNA in prokaryotes.....	46
3.4.2	Genus level comparisons.....	47
3.4.3	Phylogenetic-based comparison.....	50
3.4.4	Nucleotide composition and length of vertebrates 18S rRNA.....	54
3.5	Discussion.....	56
Chapter 4	Mutational bias affects protein evolution in flowering plants.....	61
4.1	Abstract.....	61
4.2	Introduction.....	62
4.3	Materials and Methods.....	62
4.3.1	Sources of sequence data.....	63
4.3.2	Identification and comparison of homologous sequences.....	63
4.3.3	Identifying amino acids for GC-rich and AT-rich codons.....	64

4.4 Results.....	64
4.4.1 Compositional distribution of rice and <i>Arabidopsis</i> homologous genes.....	64
4.4.2 Amino acid substitutions between rice and <i>Arabidopsis</i> homologs.....	67
4.4.3 Possible sources of compositional bias in rice genes and their encoded proteins.....	72
4.4.4 Mutational bias affects protein sequence similarity.....	75
4.5 Discussion.....	76
Chapter 5 Nucleotide content affects synonymous codon usage in rice genes.....	83
5.1 Abstract.....	83
5.2 Introduction.....	84
5.3 Materials and Methods.....	85
5.3.1 Coding sequence data.....	85
5.3.2 Identification of homologous sequences.....	86
5.3.3 Statistical analyses.....	86
5.3.3.1 G+C content and GC disparity.....	86
5.3.3.2 Codon usage indices.....	87
5.3.3.3 Correspondence analysis.....	90
5.4 Results.....	91
5.4.1 G+C content distribution and G+C disparity.....	91
5.4.2 Relative synonymous codon usage.....	92
5.4.3 Codon usage entropy.....	96
5.4.4 Effective number of codons.....	97
5.4.5 Correspondence analysis.....	100
5.4.6 Codon adaptation index.....	105
5.5 Discussion.....	107
Chapter 6 Conclusions.....	111
6.1 Thermal adaptation of rRNA genes and the genomes.....	111
6.2 Effects of nucleotide bias on codon usage and protein evolution.....	113

6.3 Future directions.....	115
6.3.1 GC content, isochores and protein evolution in vertebrates.....	115
6.3.2 Nucleotide bias, thermal adaptation and other forms of selection.....	116
6.3.3 Other perspectives.....	118
References.....	121

List of Tables

TABLE 2.1 GC content and optimal growth temperature (T_{opt} , °C) of completely sequenced genomes used in this study.....	25
TABLE 2.2 Average nucleotide composition (mean% \pm standard error) of whole genomes and 16S rRNA genes of mesophiles and thermophiles. A) Nucleotide composition of entire genome, for the 31 genomes listed in Table 1. B) Nucleotide composition of 16S rRNA genes.....	26
TABLE 2.3 Average nucleotide composition (mean% \pm standard error) of structural components of 16S rRNA genes of 44 mesophiles and thermophiles.....	27
TABLE 2.4 Correlation and regression analysis of nucleotide composition of 16S rRNA and optimal growth temperature.....	31
TABLE 2.5 The relationship between the overall G+C content of the genome and the G+C content of A) 16S rRNA unpaired regions and B) the ribosomal protein gene coding sequences.....	33
TABLE 3.1 Average GC content and sequence length of 16S rRNA stems and loops for mesophilic (< 40°C), moderately thermophilic (40-75 °C) and hyperthermophilic (\geq 75°C) bacteria (A) and archaea (B).....	46
TABLE 3.2 Average 16S rRNA GC content (%) and length (bases) of mesophilic species and thermophilic species in the same genus.....	48
TABLE 3.3 Average GC content (%) and cumulative length (bases) of stems and loops of 16S rRNA of mesophilic species and thermophilic species in the same genus.....	49
TABLE 3.4 Correlation coefficients of GCrRNA and rRNA length with optimal temperature for original data (20 species), contrasts based on the translation tree (18 species) and contrasts based on the transcription tree (20 species).....	50
TABLE 4.1 Average nucleotide contents of homologous genes in rice and <i>Arabidopsis</i> (expressed as percentages of G+C).....	66
TABLE 4.2 Exon-intron structure of rice genes and their <i>Arabidopsis</i> homologs.....	80
TABLE 5.1 A and B) Cumulative RSCU for 14005 rice and 25625 <i>Arabidopsis</i> genes. C1, C2, C3) RSCU for the three GC sets of rice genes.....	94
TABLE 5.2 Gene number and average GC content (%) of high, intermediate and low GC	