

Adaptation in genes, duplicates, families, functional modules and genomes

François Serra

October 2011

Contents

1. Introduction	2
1.1. Adaptive changes to evolutionary speed	3
1.2. Evolution, and the detection at molecular level	3
1.3. Grouping genes and finding evolutionary patterns	3
1.4. What is DNA? How genes rose?	3
1.5. Life in DNA, from genes to repetitive elements.	3
 I. Structure and dynamics of genomes	 4
2. Random-like structure of DNA	5
2.1. Introduction	6
2.2. Results and Discussion	6
2.3. Material and methods	6
2.3.1. The complexity ratio and complexity value	6
2.3.2. Complexity in strings	6
2.3.3. Simulations	6
3. Life inside genomes, dynamics and predictions	7
3.1. Genomic elements, dispersion and abundance	8
3.2. Species Abundance Diversity in genomes	8
3.3. Neutrality of SAD	8
3.4. Material Methods	8
3.4.1. Ecology	8
 II. Detection of selective pressures in genomes	 9
4. Searching for evolutionary patterns in funcionally linked group of genes	10
4.1. Introduction	11
4.2. Material and Methods	11
4.2.1. Dataset	11
Five mammals	11
6 Drosophila	11
4.2.2. Alignments	11
4.3. open on colocalization to not random	11

Contents

5. Tools, programs, methods	13
5.1. ETE-evol plugin	14
5.1.1. BRANCHED1	14
5.1.2. Protamines Rodents and Primates	14
5.2. Pipeline for study of adaptation at genomic scale	14
5.2.1. Selective pressure on duplicated genes in Drosophila	14
5.3. Phylemon	14
6. Conclusions	15

Nomenclatura

BWT Burros-Wheeler transform

CR Complexity Ratio

1. Introduction

- 1.1. Adaptive changes to evolutionary speed**
- 1.2. Evolution, and the detection at molecular level**
- 1.3. Grouping genes and finding evolutionary patterns**
- 1.4. What is DNA? How genes rose?**
- 1.5. Life in DNA, from genes to repetitive elements.**

Part I.

Structure and dynamics of genomes

2. Random-like structure of DNA

2.1. Introduction

From a biological perspective it seems obvious that DNA is something else than random mix of A, T, G or C nucleotides. Genomes are composed of functional elements as can be genes or promoters but also repetitive elements that by definition can not be random when taken together. However to what extent can we state that genomes are not a random soup of 4 letters?

This question could be solved in some sense by measuring genomes entropy. This measure presents the disadvantage that extreme cases of high entropy could correspond to *a) a specially high content of information*, entropy-based algorithms are actually used to predict or confirm automatic detection of genes [Du *et al.*2006, Gerstein *et al.*2007], *b) an exact random structure*, some work in the sense of testing the random structure of DNA have been done using entropy [Loewenstern & Yianilos1999]. However this characteristic of entropy could be only a semantic problem if we use it as a measure of relative variation in DNA complexity in genomes, and try to discern statistical patterns in the DNA sequences of different genomic element such as interspersed repeats or functional element (like protein-coding genes). This kind of description of DNA sequence complexity was already done by [Holste *et al.*2001], but only in human chromosome 22.

2.2. Results and Discussion

2.3. Material and methods

2.3.1. The complexity ratio and complexity value

Complexity Ratio CR is defined by a classical formula used in data compression [Adjero *et al.*2008], the Burros-Wheeler transform BWT [Burrows & Wheeler1994].

2.3.2. Complexity in strings

2.3.3. Simulations

3. Life inside genomes, dynamics and predictions

3. Life inside genomes, dynamics and predictions

3.1. Genomic elements, dispersion and abundance

3.2. Species Abundance Diversity in genomes

3.3. Neutrality of SAD

3.4. Material Methods

3.4.1. Ecology

Part II.

Detection of selective pressures in genomes

4. Searching for evolutionary patterns in functionally linked group of genes

4.1. Introduction

4.2. Material and Methods

4.2.1. Dataset

Five mammals

Complete genomes of 5 mammals species (*Homo sapiens*, *Pan troglodytes*, *Mus musculus*, *Rattus norvegicus* and *Canis familiaris*) were retrieved from *Ensembl* [Flicek *et al.*2011]. Also orthology prediction between each pair of species possibly done between human and the others was retrieved from *Ensembl Compara* [Vilella *et al.*2009] using biomaRt [Kinsella *et al.*2011]. Only groups of orthologs *one-to-one* with one representative of each species were kept in the final dataset.

4.1 NUMBERS

6 Drosophila

4.2.2. Alignments

Each of the group of orthologous sequences were aligned with Muscle [Edgar2004], and, once aligned sequences were cleaned with trimAl [Capella-Gutiérrez *et al.*2009] keeping all sequences but trimming alignment columns with the euristic1 method.

4.3. open on colocalization to not random

4. Searching for evolutionary patterns in functionally linked group of genes

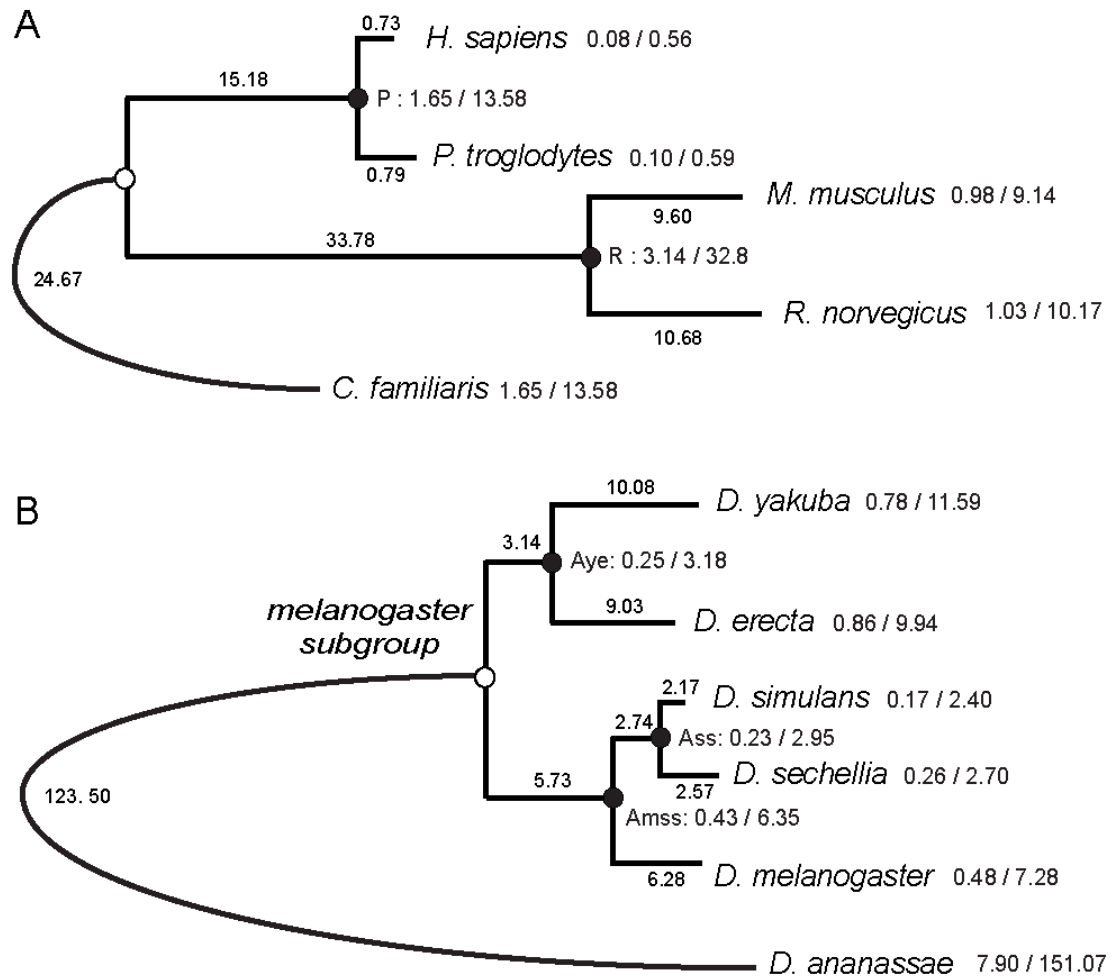


Figure 4.1.: Mammals and *Drosophila* phylogeny. blabli blob lu dkfnlskjdf

5. Tools, programs, methods

5. Tools, programs, methods

5.1. ETE-evol plugin

5.1.1. BRANCHED1

5.1.2. Protamines Rodents and Primates

5.2. Pipeline for study of adaptation at genomic scale

5.2.1. Selective pressure on duplicated genes in *Drosophila*

5.3. Phylemon

6. Conclusions

Bibliography

- [Adjeroh *et al.*2008] DONALD ADJEROH, TIM BELL, AND AMAR MUKHERJEE, The Burrows-Wheeler Transform: Data Compression, Suffix Arrays, and Pattern Matching. In *ACM SIGACT News*, vol. 41, 21–24. Springer US, Boston, MA, 2008.
- [Burrows & Wheeler1994] MICHAEL BURROWS AND DAVID J WHEELER, A block-sorting lossless data compression algorithm. *Digital SRC Research Report* **124** (1994).
- [Capella-Gutiérrez *et al.*2009] SALVADOR CAPELLA-GUTIÉRREZ, JOSÉ M SILLA-MARTÍNEZ, AND TONI GABALDÓN, trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics (Oxford, England)* **25**(15) (2009), 1972–3.
- [Du *et al.*2006] JIANG DU, JOEL S ROZOWSKY, JAN O KORBEL, ZHENG DONG D ZHANG, THOMAS E ROYCE, MARTIN H SCHULTZ, MICHAEL SNYDER, AND MARK GERSTEIN, A supervised hidden markov model framework for efficiently segmenting tiling array data in transcriptional and chIP-chip experiments: systematically incorporating validated biological knowledge. *Bioinformatics (Oxford, England)* **22**(24) (2006), 3016–24.
- [Edgar2004] ROBERT C EDGAR, MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic acids research* **32**(5) (2004), 1792–7.
- [Flicek *et al.*2011] PAUL FLICEK, M RIDWAN AMODE, DANIEL BARRELL, KATHRYN BEAL, SIMON BRENT, YUAN CHEN, PETER CLAPHAM, GUY COATES, SUSAN FAIRLEY, STEPHEN FITZGERALD, LEO GORDON, MAURICE HENDRIX, THIBAUT HOURLIER, NATHAN JOHNSON, ANDREAS KÄHÄRI, DAMIAN KEEFE, STEPHEN KEENAN, RHODA KINSELLA, FELIX KOKOCINSKI, EUGENE KULESHA, PONTUS LARSSON, IAN LONGDEN, WILLIAM McLAREN, BERT OVERDUIN, BETHAN PRITCHARD, HARPREET SINGH RIAT, DANIEL RIOS, GRAHAM R S RITCHIE, MAGALI RUFFIER, MICHAEL SCHUSTER, DANIEL SOBRAL, GIULIETTA SPUDICH, Y AMY TANG, STEPHEN TREVANION, JANA VANDROVCOVA, ALBERT J VILELLA, SIMON WHITE, STEVEN P WILDER, AMONIDA ZADISSA, JORGE ZAMORA, BRONWEN L AKEN, EWAN BIRNEY, FIONA CUNNINGHAM, IAN DUNHAM, RICHARD DURBIN, XOSÉ M FERNÁNDEZ-SUAREZ, JAVIER HERRERO, TIM J P HUBBARD, ANNE PARKER, GLENN PROCTOR, JAN VOGEL, AND STEPHEN M J SEARLE, Ensembl 2011. *Nucleic acids research* **39**(Database issue) (2011), D800–6.
- [Gerstein *et al.*2007] MARK B GERSTEIN, CAN BRUCE, JOEL S ROZOWSKY, DEYOU ZHENG, JIANG DU, JAN O KORBEL, OLOF EMANUELSSON, ZHENG DONG D ZHANG, SHERMAN WEISSMAN, AND MICHAEL SNYDER, What is a gene, post-ENCODE? History and updated definition. *Genome research* **17**(6) (2007), 669–81.
- [Holste *et al.*2001] DIRK HOLSTE, IVO GROSSE, AND HANSPETER HERZEL, Statistical analysis of the DNA sequence of human chromosome 22. *Physical Review E* **64**(4) (2001), 1–9.
- [Kinsella *et al.*2011] RHODA J. KINSELLA, ANDREAS KÄHÄRI, SYED HAIDER, JORGE ZAMORA, GLENN PROCTOR, GIULIETTA SPUDICH, JEFF ALMEIDA-KING, DANIEL STAINES, PAUL

Bibliography

- DERWENT, ARNAUD KERHORNOU, PAUL KERSEY, AND PAUL FLICEK, Ensembl BioMarts: a hub for data retrieval across taxonomic space. *Database : the journal of biological databases and curation* **2011** (2011), p. bar030.
- [Loewenstern & Yianilos1999] DAVID M LOEWENSTERN AND PETER N YIANILOS, Significantly lower entropy estimates for natural DNA sequences. *Journal of computational biology : a journal of computational molecular cell biology* **6**(1) (1999), 125–42.
- [Vilella *et al.*2009] ALBERT J VILELLA, JESSICA SEVERIN, ABEL URETA-VIDAL, LI HENG, RICHARD DURBIN, AND EWAN BIRNEY, EnsemblCompara GeneTrees: Complete, duplication-aware phylogenetic trees in vertebrates. *Genome research* **19**(2) (2009), 327–35.

List of Figures

4.1. Mammals and <i>Drosophila</i> phylogeny	12
--	----

List of Tables