

# Informational, Ecological and System Approaches for Complete Genome Analysis

François Serra<sup>1</sup>

Supervised by: Hernán Dopazo<sup>1</sup>

<sup>1</sup>Evolutionary Genomic Laboratory,  
Centro de Investigación Príncipe Felipe

12<sup>th</sup> December 2012



# Index

- 1 Introduction
  - Evolution only makes sense in light of neutrality
  - Overview
- 2 Random-like structure of DNA
  - Background
  - Methodology
  - Results
- 3 Ecology of genetic elements
  - Background
  - Results
- 4 Genomic study of selective pressures in group of genes
  - Background
  - Methodology
  - Results

# Index

## 1 Introduction

- Evolution only makes sense in light of neutrality
- Overview

## 2 Random-like structure of DNA

- Background
- Methodology
- Results

## 3 Ecology of genetic elements

- Background
- Results

## 4 Genomic study of selective pressures in group of genes

- Background
- Methodology
- Results

# Classification of evolutionary changes

- Evolutionary changes affect the inherited characteristics of a biological population
- **Changes are not directional**, evolution may accept or reject them through the elimination of individuals
- According to the influence of natural selection, we can classify **changes** in three categories:
  - **advantageous**: individuals carrying them may improve their fitness and thus their chances to produce offspring
  - **deleterious**: result in a loss of fitness, and may be removed from populations
  - **neutral**: they represent the narrow area that segregates advantageous and deleterious changes, they may not affect the fitness, and thus, may escape natural selection

# Classification of evolutionary changes

- Evolutionary changes affect the inherited characteristics of a biological population
- **Changes are not directional**, evolution may accept or reject them through the elimination of individuals
- According to the influence of natural selection, we can classify **changes** in three categories:
  - **advantageous**: individuals carrying them may improve their fitness and thus their chances to produce offspring
  - **deleterious**: result in a loss of fitness, and may be removed from populations
  - **neutral**: they represent the narrow area that segregates advantageous and deleterious changes, they may not affect the fitness, and thus, may escape natural selection

# Classification of evolutionary changes

- Evolutionary changes affect the inherited characteristics of a biological population
- **Changes are not directional**, evolution may accept or reject them through the elimination of individuals
- According to the influence of natural selection, we can classify **changes** in three categories:
  - **advantageous**: individuals carrying them may improve their fitness and thus their chances to produce offspring
  - **deleterious**: result in a loss of fitness, and may be removed from populations
  - **neutral**: they represent the narrow area that segregates advantageous and deleterious changes, they may not affect the fitness, and thus, may escape natural selection

# Classification of evolutionary changes

- Evolutionary changes affect the inherited characteristics of a biological population
- **Changes are not directional**, evolution may accept or reject them through the elimination of individuals
- According to the influence of natural selection, we can classify **changes** in three categories:
  - **advantageous**: individuals carrying them may improve their fitness and thus their chances to produce offspring
  - **deleterious**: result in a loss of fitness, and may be removed from populations
  - **neutral**: they represent the narrow area that segregates advantageous and deleterious changes, they may not affect the fitness, and thus, may escape natural selection

# Classification of evolutionary changes

- Evolutionary changes affect the inherited characteristics of a biological population
- **Changes are not directional**, evolution may accept or reject them through the elimination of individuals
- According to the influence of natural selection, we can classify **changes** in three categories:
  - **advantageous**: individuals carrying them may improve their fitness and thus their chances to produce offspring
  - **deleterious**: result in a loss of fitness, and may be removed from populations
  - **neutral**: they represent the narrow area that segregates advantageous and deleterious changes, they may not affect the fitness, and thus, may escape natural selection

# Classification of evolutionary changes

- Evolutionary changes affect the inherited characteristics of a biological population
- **Changes are not directional**, evolution may accept or reject them through the elimination of individuals
- According to the influence of natural selection, we can classify **changes** in three categories:
  - **advantageous**: individuals carrying them may improve their fitness and thus their chances to produce offspring
  - **deleterious**: result in a loss of fitness, and may be removed from populations
  - **neutral**: they represent the narrow area that segregates advantageous and deleterious changes, they may not affect the fitness, and thus, may escape natural selection

# Importance of neutral changes in molecular evolution

- Since their definition by Charles Darwin, neutral changes were thought to be **minor actors** in the process of evolution.
- First molecular data and the estimation of mutation rate much higher than expected
  - Jack Lester King and Thomas H. Jukes, and independently Motoo Kimura to propose the *neutral theory of molecular evolution*.

# Importance of neutral changes in molecular evolution

- Since their definition by Charles Darwin, neutral changes were thought to be **minor actors** in the process of evolution.
- First molecular data and the estimation of mutation rate much higher than expected
  - Jack Lester King and Thomas H. Jukes, and independently Motoo Kimura to propose the *neutral theory of molecular evolution*.

# Importance of neutral changes in molecular evolution

- Since their definition by Charles Darwin, neutral changes were thought to be **minor actors** in the process of evolution.
- First molecular data and the estimation of mutation rate much higher than expected
  - Jack Lester King and Thomas H. Jukes, and independently Motoo Kimura to propose the *neutral theory of molecular evolution*.

# Neutrality in Ecology

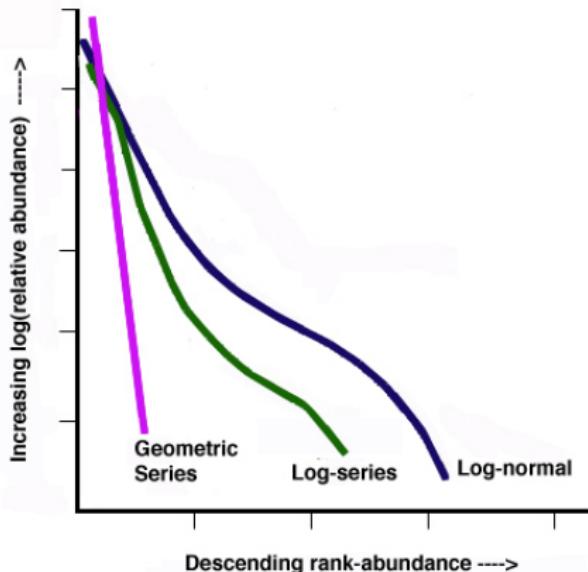
neutral distribution of species' abundances in ecosystems

- Models used were based on the concept of ecological niches. Differences observed between species were thought to be mostly adaptive
- Statistical models able to reproduce the pattern of distribution of species' abundances in ecosystems (Motomura, Fisher, Preston)
- **neutral models** had to wait for the late 20<sup>th</sup> century (Hubbell)

# Neutrality in Ecology

neutral distribution of species' abundances in ecosystems

- Models used were based on the concept of ecological niches. Differences observed between species were thought to be mostly adaptive
- Statistical models able to reproduce the pattern of distribution of species' abundances in ecosystems (Motomura, Fisher, Preston)
- neutral models had to wait for the late 20<sup>th</sup> century (Hubbell)

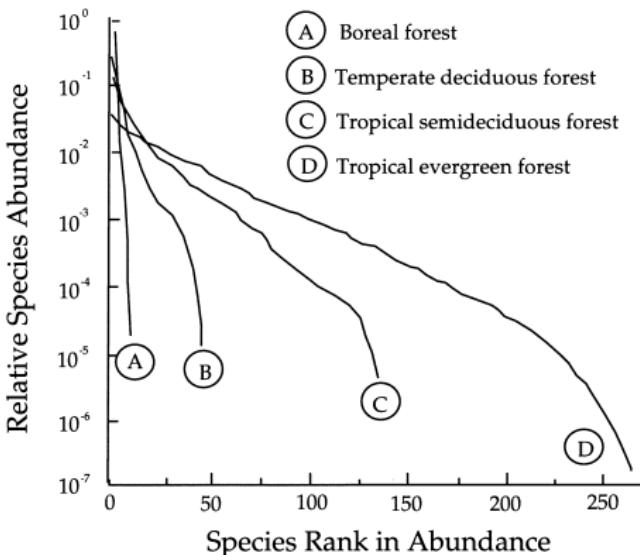


adapted from [Magurran - (1988)]

# Neutrality in Ecology

neutral distribution of species' abundances in ecosystems

- Models used were based on the concept of ecological niches. Differences observed between species were thought to be mostly adaptive
- Statistical models able to reproduce the pattern of distribution of species' abundances in ecosystems (Motomura, Fisher, Preston)
- neutral models had to wait for the late 20<sup>th</sup> century (Hubbell)

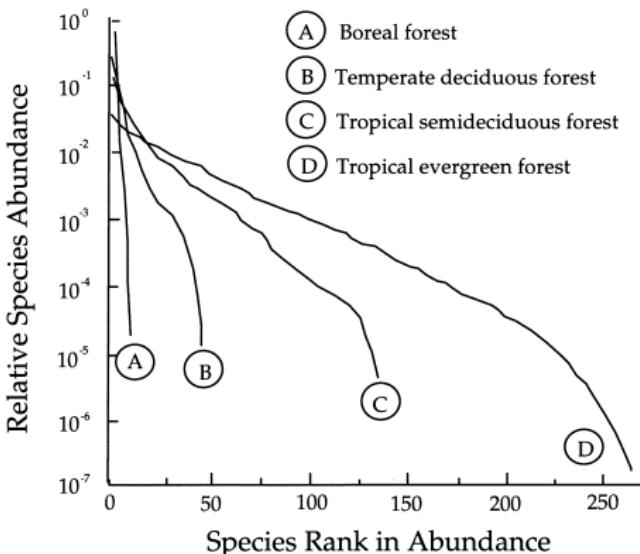


[Hubbell - (2001)]

# Neutrality in Ecology

neutral distribution of species' abundances in ecosystems

- Models used were based on the concept of ecological niches. Differences observed between species were thought to be mostly adaptive
- Statistical models able to reproduce the pattern of distribution of species' abundances in ecosystems (Motomura, Fisher, Preston)
- **neutral models** had to wait for the late 20<sup>th</sup> century (Hubbell)

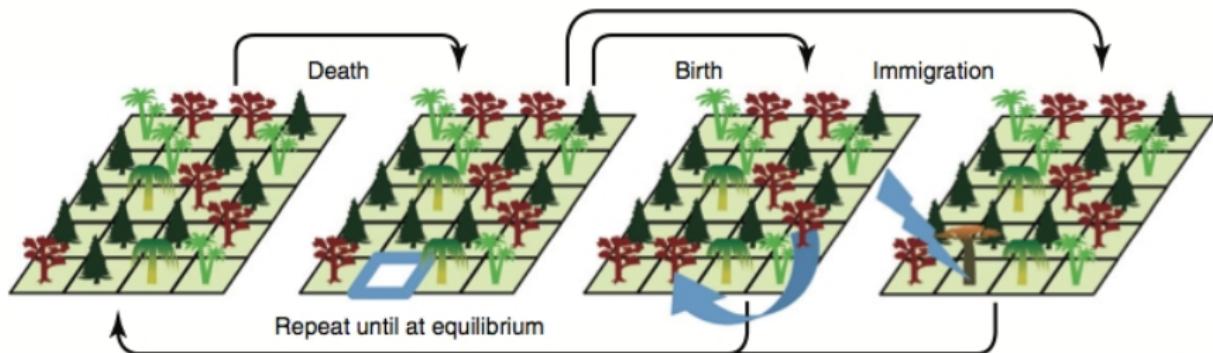


[Hubbell - (2001)]

# Neutrality in Ecology

## Unified Neutral Theory of Biodiversity and Biogeography (UNTB)

Assumes that diversity in a local community of individuals is maintained by migration from the metacommunity at a constant rate ( $m$ ). Births and deaths in the local community occur at constant rates per generation, **regardless of the species**



TRENDS in Ecology & Evolution

[Rosindell et al. - Trends in Ecology Evolution (2011)]

# Index

## 1 Introduction

- Evolution only makes sense in light of neutrality
- **Overview**

## 2 Random-like structure of DNA

- Background
- Methodology
- Results

## 3 Ecology of genetic elements

- Background
- Results

## 4 Genomic study of selective pressures in group of genes

- Background
- Methodology
- Results

# Objectives – Overview

The work presented here consists in three different approximations in the analysis of genomes, each expressing a neutral model and its corresponding deviations.

- **Informational:** where genomes are viewed as simple sequences with a given informational content. We will define a common structure of genomes and show how they are characterized by a universal quasi-random structure.
- **Ecological:** here a genetic elements populating genomes are assimilated to biological species in ecosystems. We will apply ecological models on genomes in order to expose the neutral pattern behind their composition.
- **System:** focusing on proteins working together to complete a function we present a new strategy to find natural selection fingerprint at this level of study.

# Objectives – Overview

The work presented here consists in three different approximations in the analysis of genomes, each expressing a neutral model and its corresponding deviations.

- **Informational:** where genomes are viewed as simple sequences with a given informational content. We will define a common structure of genomes and show how they are characterized by a universal quasi-random structure.
- **Ecological:** here a genetic elements populating genomes are assimilated to biological species in ecosystems. We will apply ecological models on genomes in order to expose the neutral pattern behind their composition.
- **System:** focusing on proteins working together to complete a function we present a new strategy to find natural selection fingerprint at this level of study.

# Objectives – Overview

The work presented here consists in three different approximations in the analysis of genomes, each expressing a neutral model and its corresponding deviations.

- **Informational:** where genomes are viewed as simple sequences with a given informational content. We will define a common structure of genomes and show how they are characterized by a universal quasi-random structure.
- **Ecological:** here a genetic elements populating genomes are assimilated to biological species in ecosystems. We will apply ecological models on genomes in order to expose the neutral pattern behind their composition.
- **System:** focusing on proteins working together to complete a function we present a new strategy to find natural selection fingerprint at this level of study.

# Objectives – Overview

The work presented here consists in three different approximations in the analysis of genomes, each expressing a neutral model and its corresponding deviations.

- **Informational:** where genomes are viewed as simple sequences with a given informational content. We will define a common structure of genomes and show how they are characterized by a universal quasi-random structure.
- **Ecological:** here a genetic elements populating genomes are assimilated to biological species in ecosystems. We will apply ecological models on genomes in order to expose the neutral pattern behind their composition.
- **System:** focusing on proteins working together to complete a function we present a new strategy to find natural selection fingerprint at this level of study.

# Index

## 1 Introduction

- Evolution only makes sense in light of neutrality
- Overview

## 2 Random-like structure of DNA

- Background
- Methodology
- Results

## 3 Ecology of genetic elements

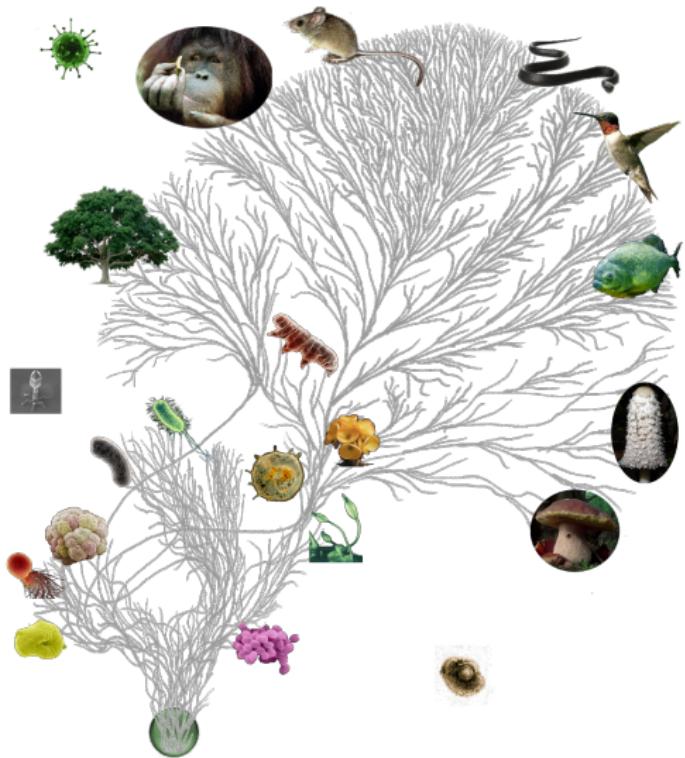
- Background
- Results

## 4 Genomic study of selective pressures in group of genes

- Background
- Methodology
- Results

# Adaptive raise in organism complexity

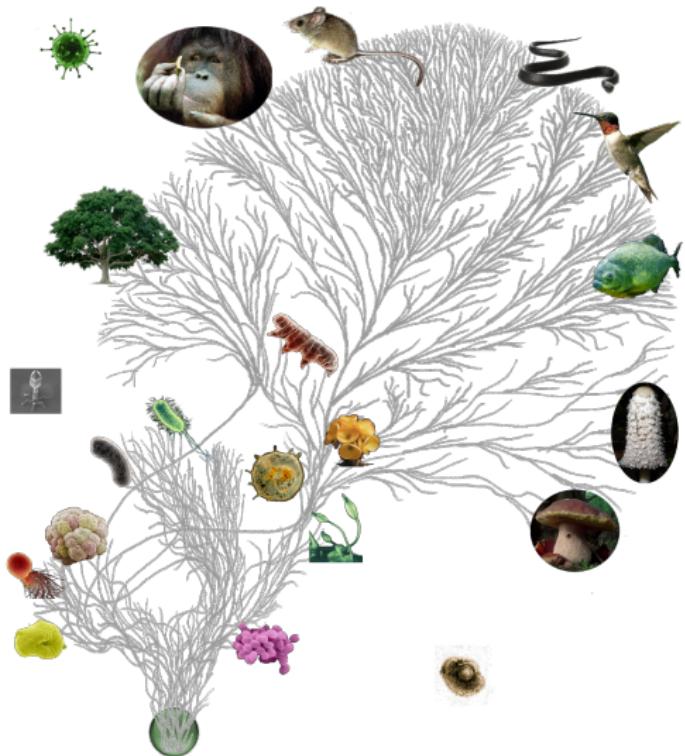
- Darwinian evolution does not encompass directional change nor global adaptive advance
- Not easy to define biological complexity
- At genomic level we would expect that the quantity of hereditary information is proportional to the level of complexity



adapted from [Maddison & Schulz - (2007)]

# Adaptive raise in organism complexity

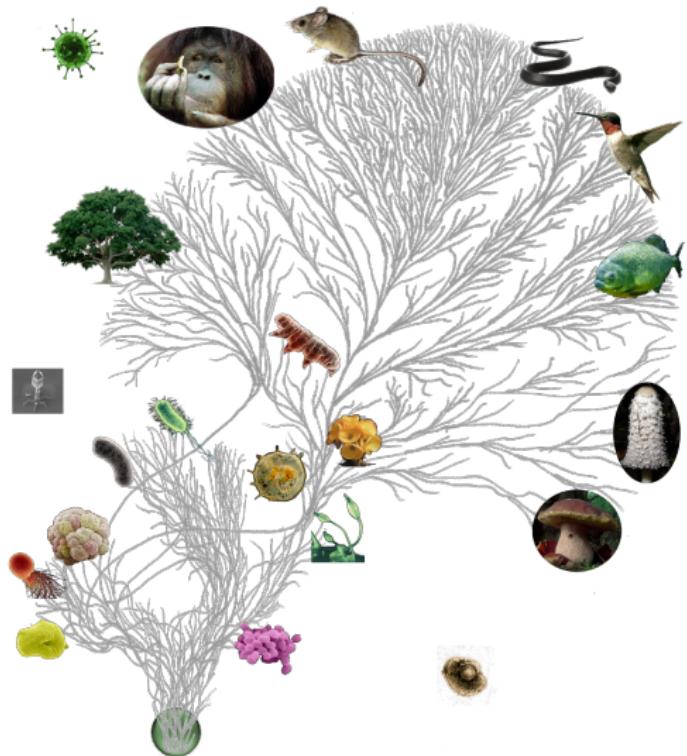
- Darwinian evolution does not encompass directional change nor global adaptive advance
- Not easy to define biological complexity
- At genomic level we would expect that the quantity of hereditary information is proportional to the level of complexity



adapted from [Maddison & Schulz - (2007)]

# Adaptive raise in organism complexity

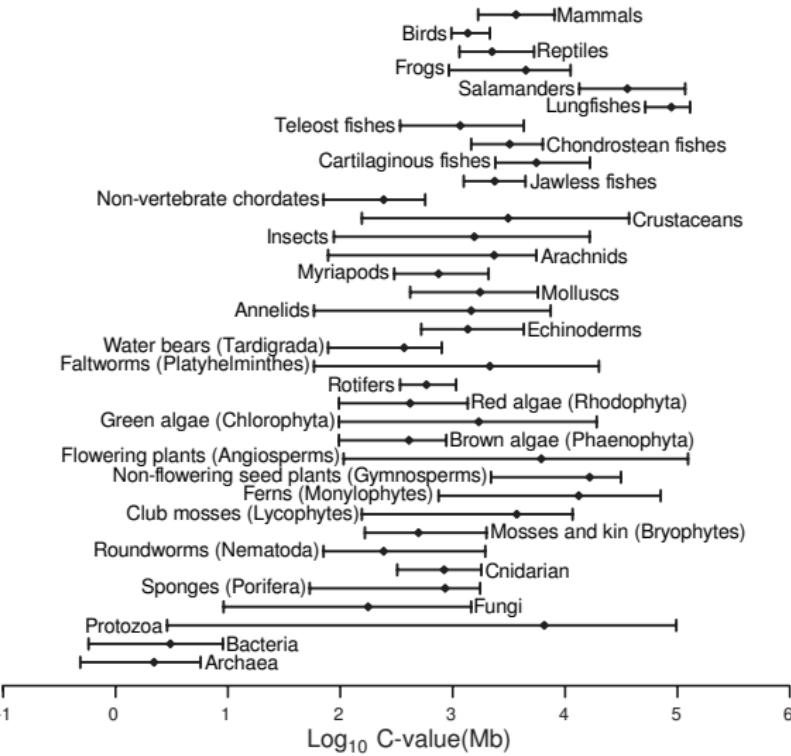
- Darwinian evolution does not encompass directional change nor global adaptive advance
- Not easy to define biological complexity
- At genomic level we would expect that the quantity of hereditary information is proportional to the level of complexity



adapted from [Maddison & Schulz - (2007)]

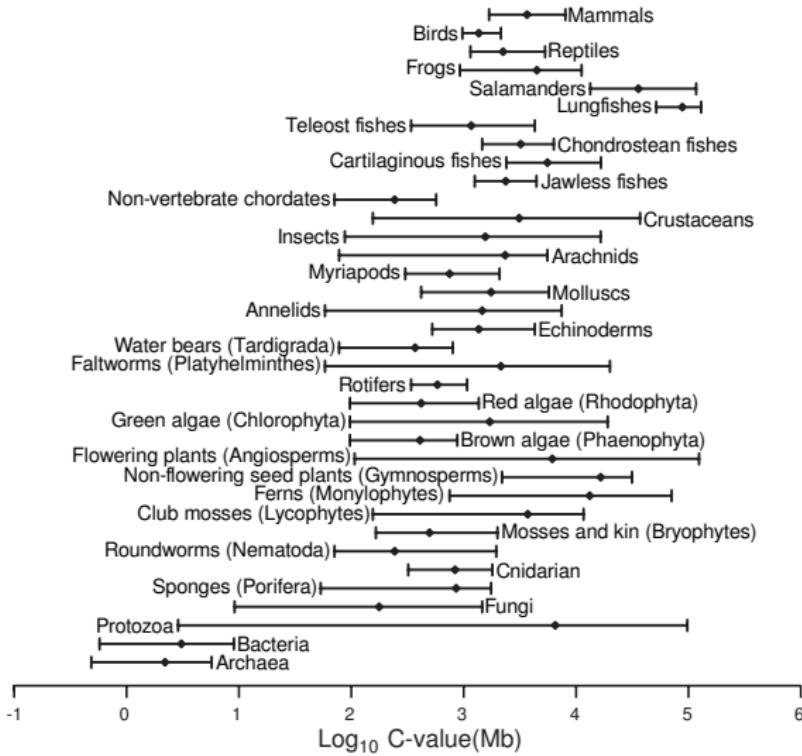
# Challenging the C-value paradox

- DNA content (C-value) do not correlate with biological complexity
- What can be summarized as:
  - the search of patterns in the informational content of genomes
  - the dynamics underlying the distribution and appearance of non-coding DNA



# Challenging the C-value paradox

- DNA content (C-value) do not correlate with biological complexity
- What can be summarized as:
  - the search of patterns in the informational content of genomes
  - the dynamics underlying the distribution and appearance of non-coding DNA



# Informational content

- In order to normalize genomes DNA content, instead of working with direct measures of genome size, we work with informational content.
- Informational content corresponds to DNA complexity and may be related to the compressed genome size
- Indeed the complexity value we use is based on classical algorithms of data compression

# Index

## 1 Introduction

- Evolution only makes sense in light of neutrality
- Overview

## 2 Random-like structure of DNA

- Background
- **Methodology**
- Results

## 3 Ecology of genetic elements

- Background
- Results

## 4 Genomic study of selective pressures in group of genes

- Background
- Methodology
- Results

# The Genome/Chromosome Complexity Value

Given a sequence,  $seq = AACCTTCGTAGCATGG$ :

#	Rotating sequence	<i>I.</i>	BWT	Char. list	MTF
0	AACCTTCGTAGCATGG	0	G	G a t c	0
1	ACCTTCGTAGCATGG A	1	A	G A t c	1
2	CCTTCGTAGCATGG AA	5	T	A g T c	2
3	TTCTCGTAGCATGG AAC	7	C	T a g C	3
4	TCGTAGCATGG AACCC	15	G	C t a G	3
5	CGTAGCATGG AACCT	13	A	G c t A	3
6	GTAGCATGG AACCTT	6	T	A g c T	3
7	TAGCATGG AACCTTC	11	C	T a g C	3
8	AGCATGG AACCTTCG	12	G	C t a G	3
9	GCATGG AACCTTCGT	2	A	G c t A	3
10	CATGG AACCTTCGT	9	T	A g c T	3
11	ATGG AACCTTCGTAG	4	C	T a g C	3
12	TGG AACCTTCGTAGC	3	G	C t a G	3
13	GG AACCTTCGTAGCA	14	T	G c T a	2
14	G AACCTTCGTAGCAT	10	A	T g c A	3
15	A AACCTTCGTAGCATG	8	C	A t g C	3

$$\diamond CR(seq) = E(MTF(BWT(seq))) = E(0, 1, 2, 3, 3, 3, 3, 3, 3, 3, 3, 3, 2, 3, 3) = 0.593$$

- 54 species (virus to mammals) with genomes from 1.6Kb to 3.4Gb

# The Genome/Chromosome Complexity Value

Given a sequence,  $seq = AACCTTCGTAGCATGG$ :

#	Rotating sequence	<i>I.</i>	BWT	Char. list	MTF
0	AACCTTCGTAGCATGG	0	G	G a t c	0
1	ACCTTCGTAGCATGG A	1	A	G A t c	1
2	CCTTCGTAGCATGG AA	5	T	A g T c	2
3	TTCTCGTAGCATGG AAC	7	C	T a g C	3
4	TCGTAGCATGG AACCC	15	G	C t a G	3
5	CGTAGCATGG AACCT	13	A	G c t A	3
6	GTAGCATGG AACCTT	6	T	A g c T	3
7	TAGCATGG AACCTTC	11	C	T a g C	3
8	AGCATGG AACCTTCG	12	G	C t a G	3
9	GCATGG AACCTTCGT	2	A	G c t A	3
10	CATGG AACCTTCGT	9	T	A g c T	3
11	ATGG AACCTTCGTAG	4	C	T a g C	3
12	TGG AACCTTCGTAGC	3	G	C t a G	3
13	GG AACCTTCGTAGCA	14	T	G c T a	2
14	G AACCTTCGTAGCAT	10	A	T g c A	3
15	A AACCTTCGTAGCATG	8	C	A t g C	3

$$\Leftrightarrow CR(seq) = E(MTF(BWT(seq))) = E(0, 1, 2, 3, 3, 3, 3, 3, 3, 3, 3, 3, 2, 3, 3) = 0.593$$

- 54 species (virus to mammals) with genomes from 1.6Kb to 3.4Gb

# Index

## 1 Introduction

- Evolution only makes sense in light of neutrality
- Overview

## 2 Random-like structure of DNA

- Background
- Methodology
- **Results**

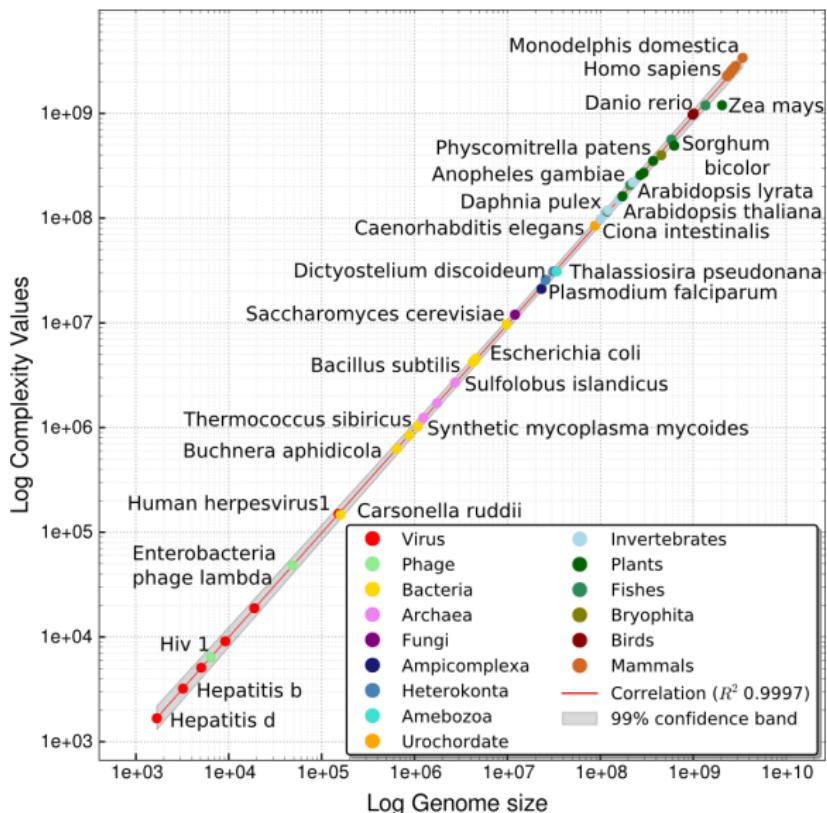
## 3 Ecology of genetic elements

- Background
- Results

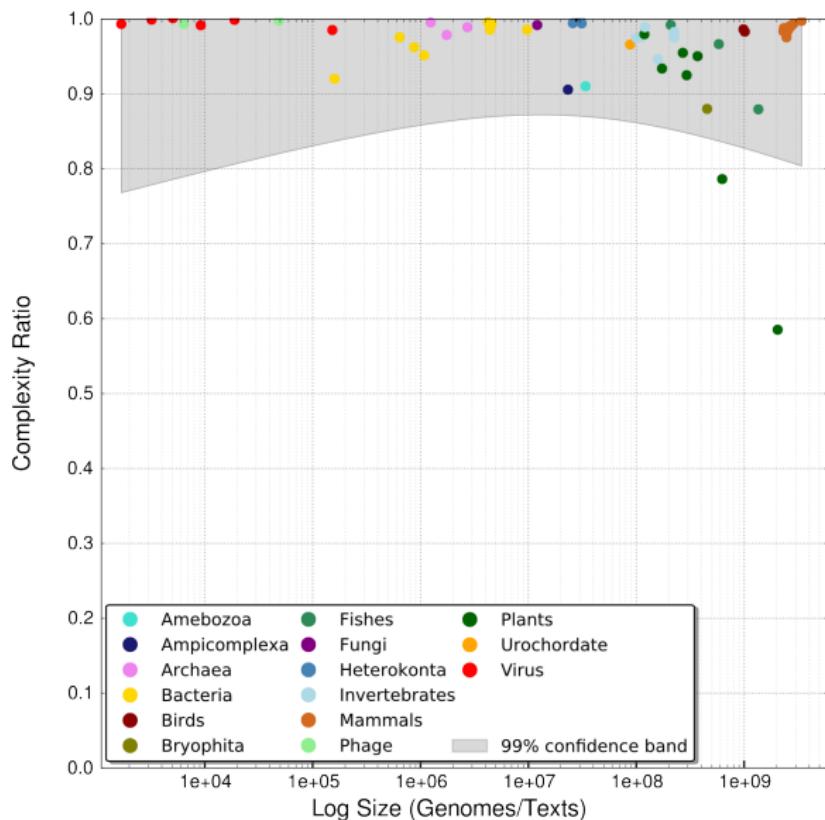
## 4 Genomic study of selective pressures in group of genes

- Background
- Methodology
- Results

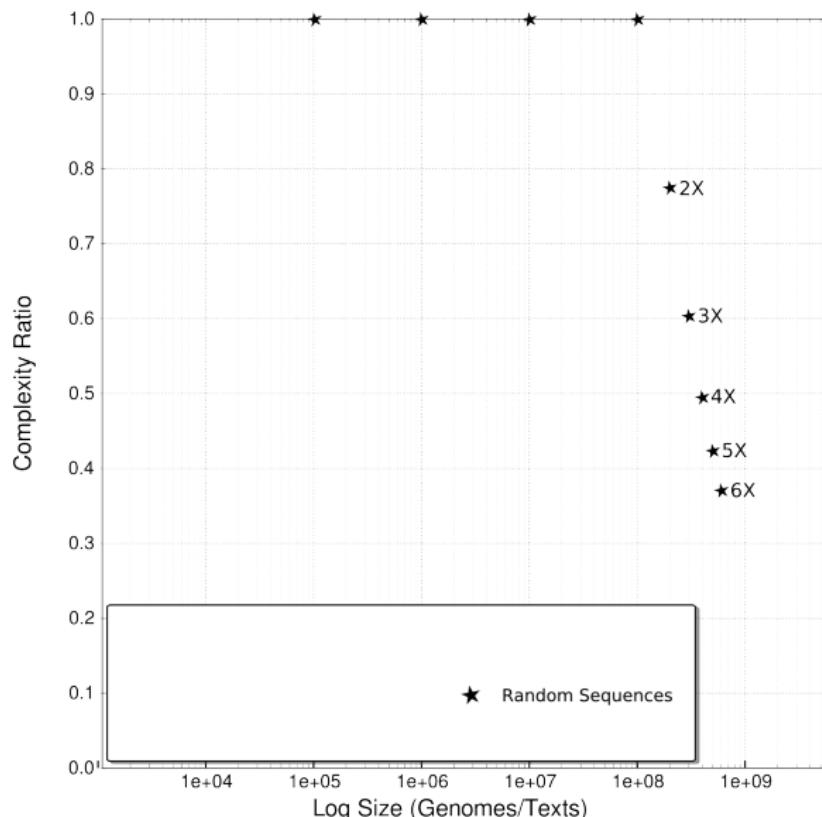
# Complexity Value of genomes



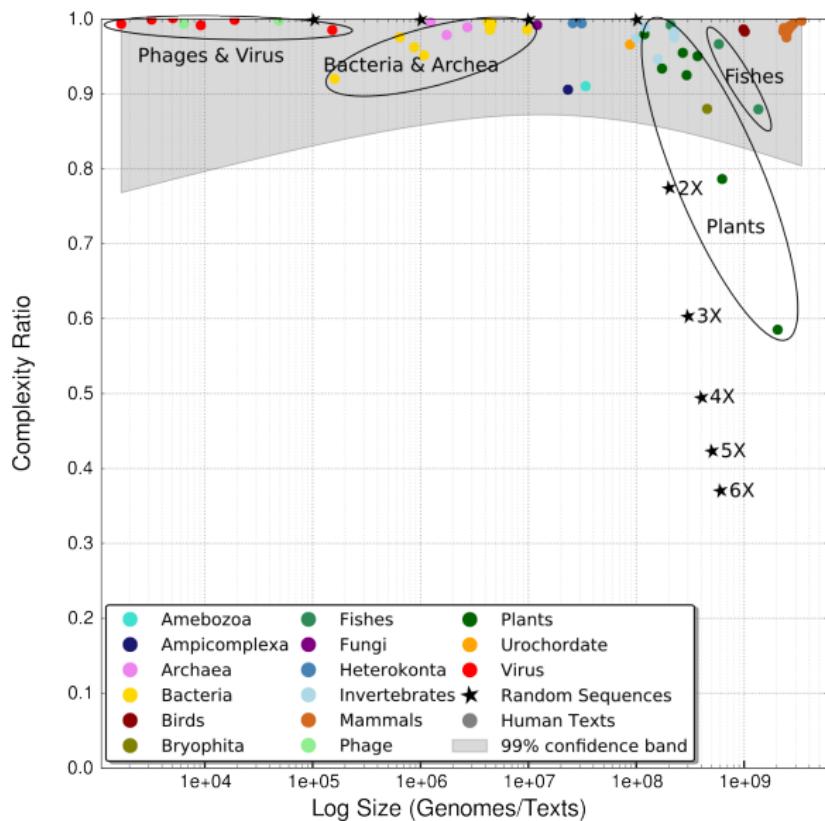
# Complexity Ratio of genomes



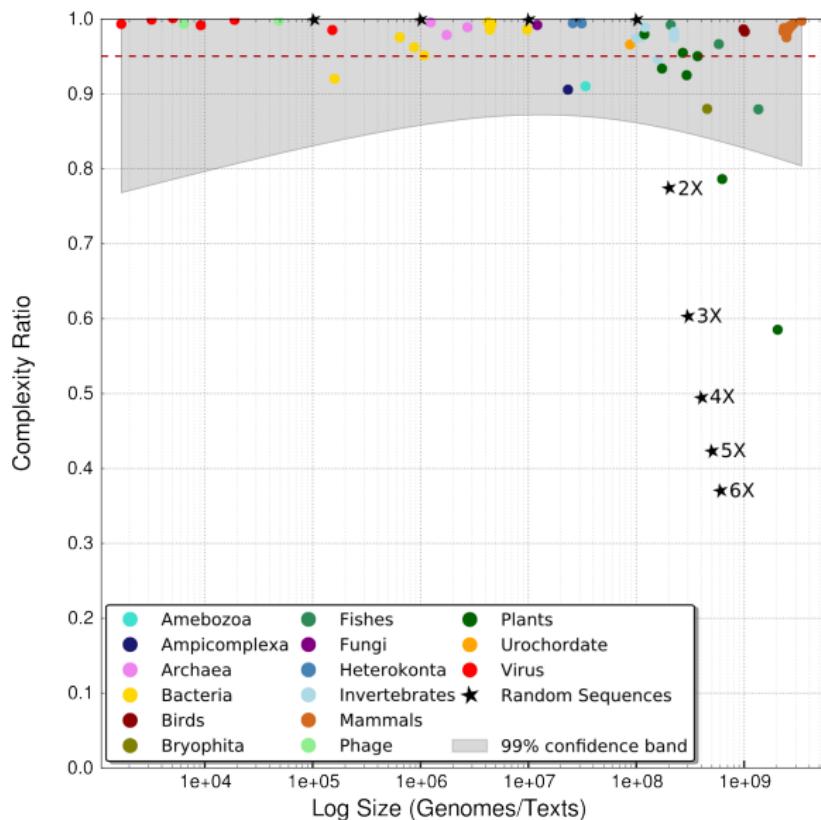
# Complexity Ratio of genomes



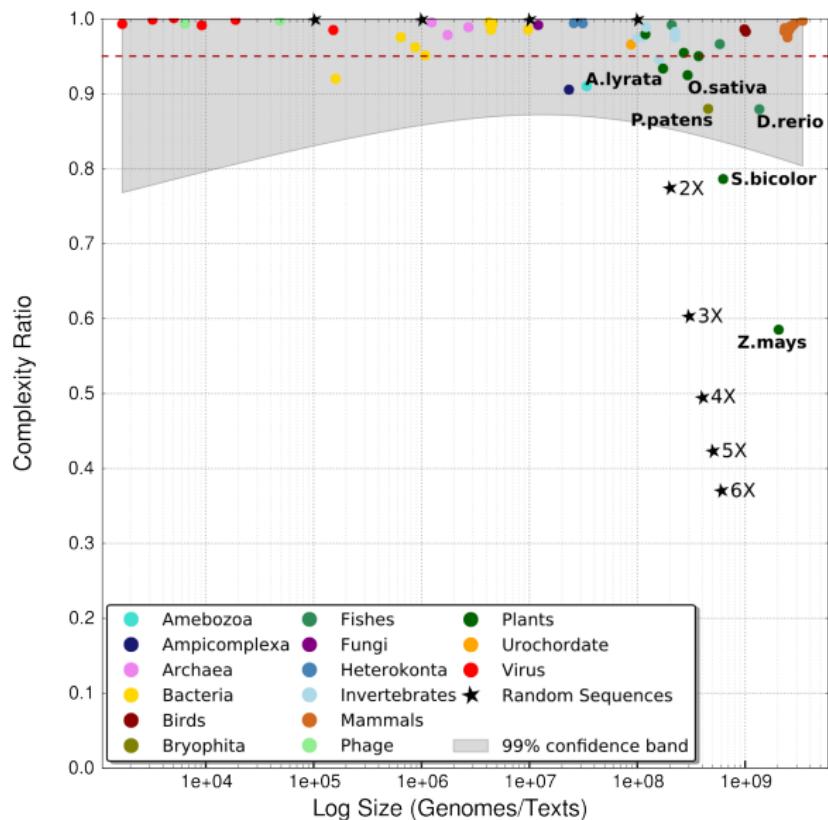
# Complexity Ratio of genomes



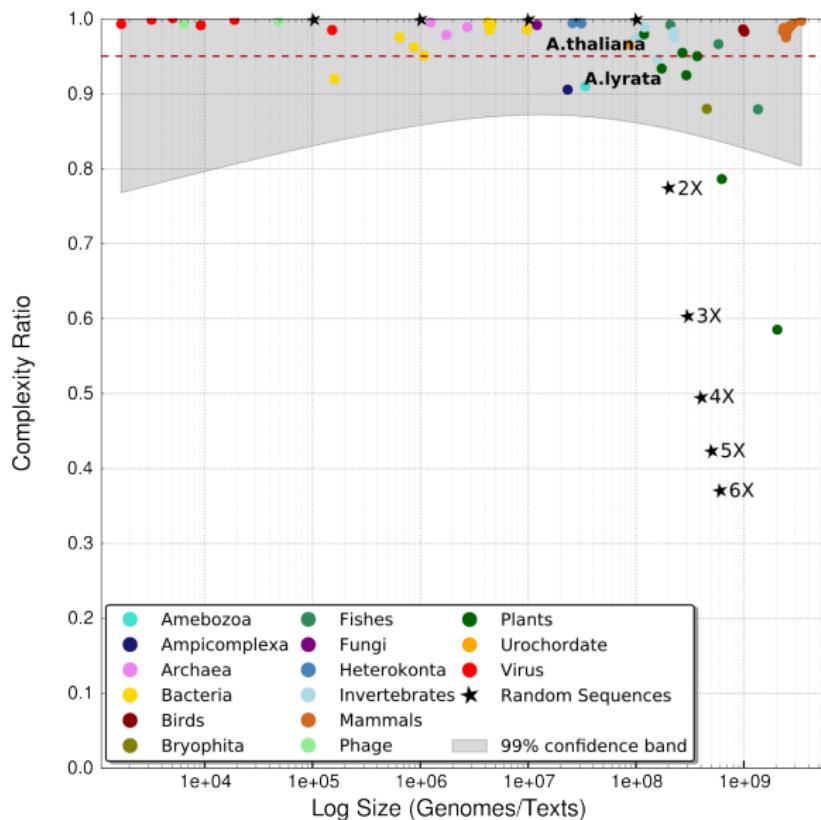
# Complexity Ratio of genomes



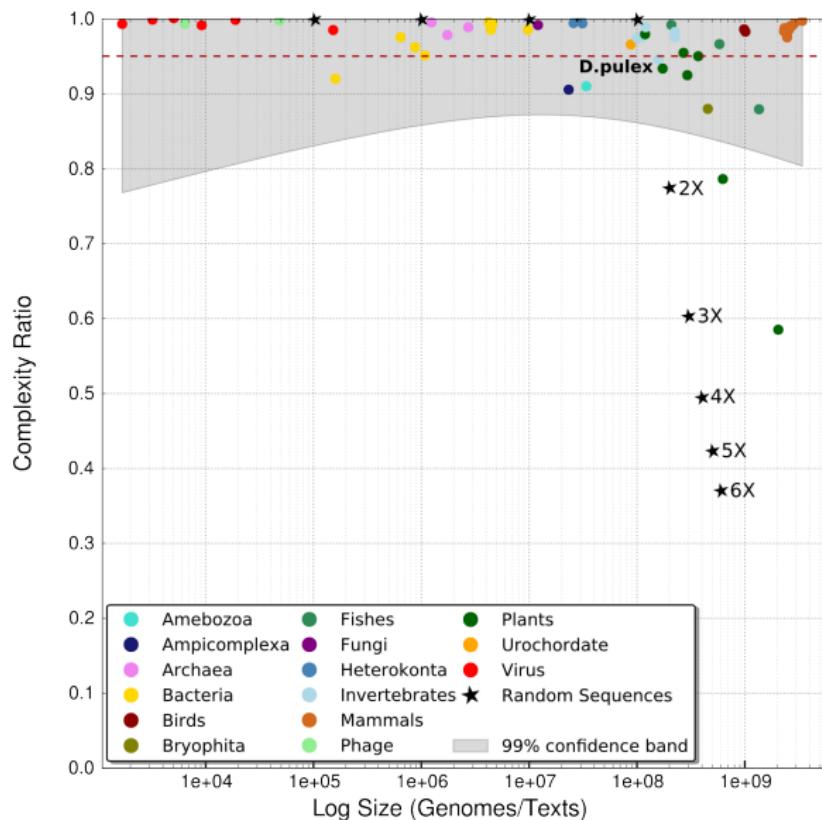
# Complexity Ratio of genomes



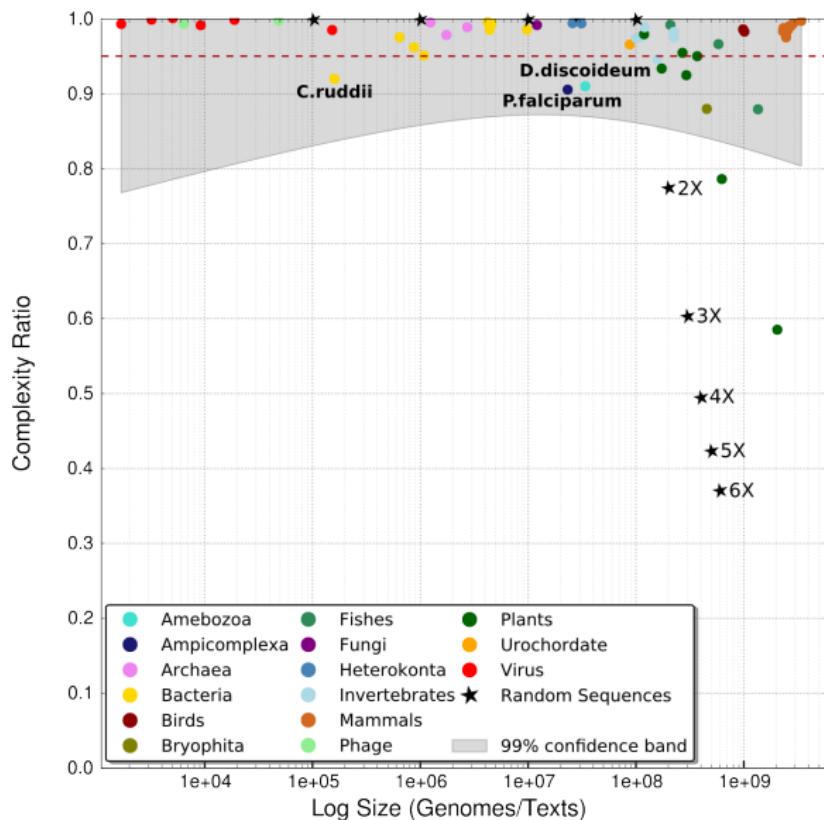
# Complexity Ratio of genomes



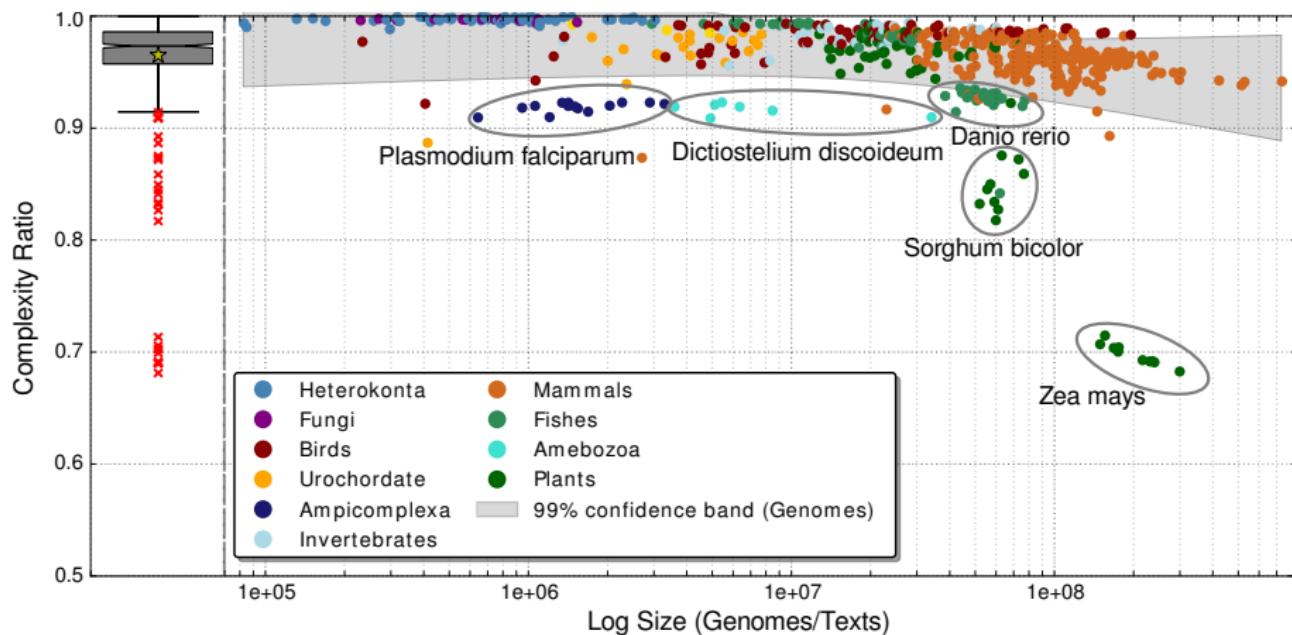
# Complexity Ratio of genomes



# Complexity Ratio of genomes



# Complexity Ratio of chromosomes



# Hypotheses

- These observations led us to hypothesize that:
  - A quasi-random combinatorial structure of DNA is a universal feature of non-polypliod genomes
  - The fate of polypliod genomes is to reach almost maximal CR. Increases in CR occur as a function of time
  - Genome complexity only increases through rounds of duplications followed by the divergence of duplicates during evolution
- These hypotheses can be nullified in some specific cases:
  - Genomes of recent polypliod species evidencing high CR, e.g.: species that have undergone significant genome reduction
  - Genomes evidencing a non-random DNA structure (low CR), due to a very strong GC content bias

# Hypotheses

- These observations led us to hypothesize that:
  - A quasi-random combinatorial structure of DNA is a universal feature of non-polypliod genomes
  - The fate of polypliod genomes is to reach almost maximal CR. Increases in CR occur as a function of time
  - Genome complexity only increases through rounds of duplications followed by the divergence of duplicates during evolution
- These hypotheses can be nullified in some specific cases:
  - Genomes of recent polypliod species evidencing high CR, e.g.: species that have undergone significant genome reduction
  - Genomes evidencing a non-random DNA structure (low CR), due to a very strong GC content bias

# Hypotheses

- These observations led us to hypothesize that:
  - A quasi-random combinatorial structure of DNA is a universal feature of non-polypliod genomes
  - The fate of polypliod genomes is to reach almost maximal CR. Increases in CR occur as a function of time
  - Genome complexity only increases through rounds of duplications followed by the divergence of duplicates during evolution
- These hypotheses can be nullified in some specific cases:
  - Genomes of recent polypliod species evidencing high CR, e.g.: species that have undergone significant genome reduction
  - Genomes evidencing a non-random DNA structure (low CR), due to a very strong GC content bias

# Hypotheses

- These observations led us to hypothesize that:
  - A quasi-random combinatorial structure of DNA is a universal feature of non-polypliod genomes
  - The fate of polypliod genomes is to reach almost maximal CR. Increases in CR occur as a function of time
  - Genome complexity only increases through rounds of duplications followed by the divergence of duplicates during evolution
- These hypotheses can be nullified in some specific cases:
  - Genomes of recent polypliod species evidencing high CR, e.g.: species that have undergone significant genome reduction
  - Genomes evidencing a non-random DNA structure (low CR), due to a very strong GC content bias

# Hypotheses

- These observations led us to hypothesize that:
  - A quasi-random combinatorial structure of DNA is a universal feature of non-polypliod genomes
  - The fate of polypliod genomes is to reach almost maximal CR. Increases in CR occur as a function of time
  - Genome complexity only increases through rounds of duplications followed by the divergence of duplicates during evolution
- These hypotheses can be nullified in some specific cases:
  - Genomes of recent polypliod species evidencing high CR, e.g.: species that have undergone significant genome reduction
  - Genomes evidencing a non-random DNA structure (low CR), due to a very strong GC content bias

# Hypotheses

- These observations led us to hypothesize that:
  - A quasi-random combinatorial structure of DNA is a universal feature of non-polypliod genomes
  - The fate of polypliod genomes is to reach almost maximal CR. Increases in CR occur as a function of time
  - Genome complexity only increases through rounds of duplications followed by the divergence of duplicates during evolution
- These hypotheses can be nullified in some specific cases:
  - Genomes of recent polypliod species evidencing high CR, e.g.: species that have undergone significant genome reduction
  - Genomes evidencing a non-random DNA structure (low CR), due to a very strong GC content bias

# Hypotheses

- These observations led us to hypothesize that:
  - A quasi-random combinatorial structure of DNA is a universal feature of non-polypliod genomes
  - The fate of polypliod genomes is to reach almost maximal CR. Increases in CR occur as a function of time
  - Genome complexity only increases through rounds of duplications followed by the divergence of duplicates during evolution
- These hypotheses can be nullified in some specific cases:
  - Genomes of recent polypliod species evidencing high CR, e.g.: species that have undergone significant genome reduction
  - Genomes evidencing a non-random DNA structure (low CR), due to a very strong GC content bias

# Index

## 1 Introduction

- Evolution only makes sense in light of neutrality
- Overview

## 2 Random-like structure of DNA

- Background
- Methodology
- Results

## 3 Ecology of genetic elements

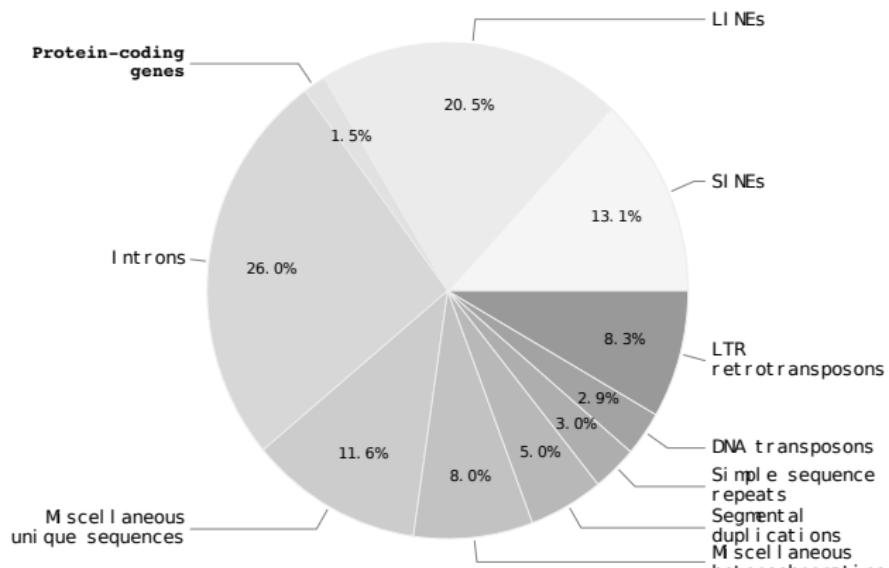
- **Background**
- Results

## 4 Genomic study of selective pressures in group of genes

- Background
- Methodology
- Results

# Variation of families of genetic elements in eukaryote

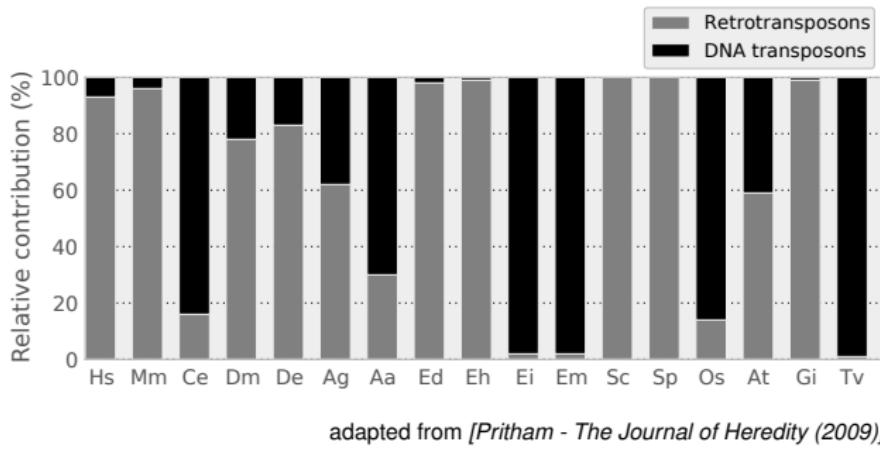
- Genomes are populated with genetic elements, coding or non-coding
- Proportions of these families of genetic elements vary greatly between species



adapted from [Lander et al. - Nature (2001)]

# Variation of families of genetic elements in eukaryote

- Genomes are populated with genetic elements, coding or non-coding
- Proportions of these families of genetic elements vary greatly between species



adapted from [Pritham - *The Journal of Heredity* (2009)]

# Index

## 1 Introduction

- Evolution only makes sense in light of neutrality
- Overview

## 2 Random-like structure of DNA

- Background
- Methodology
- Results

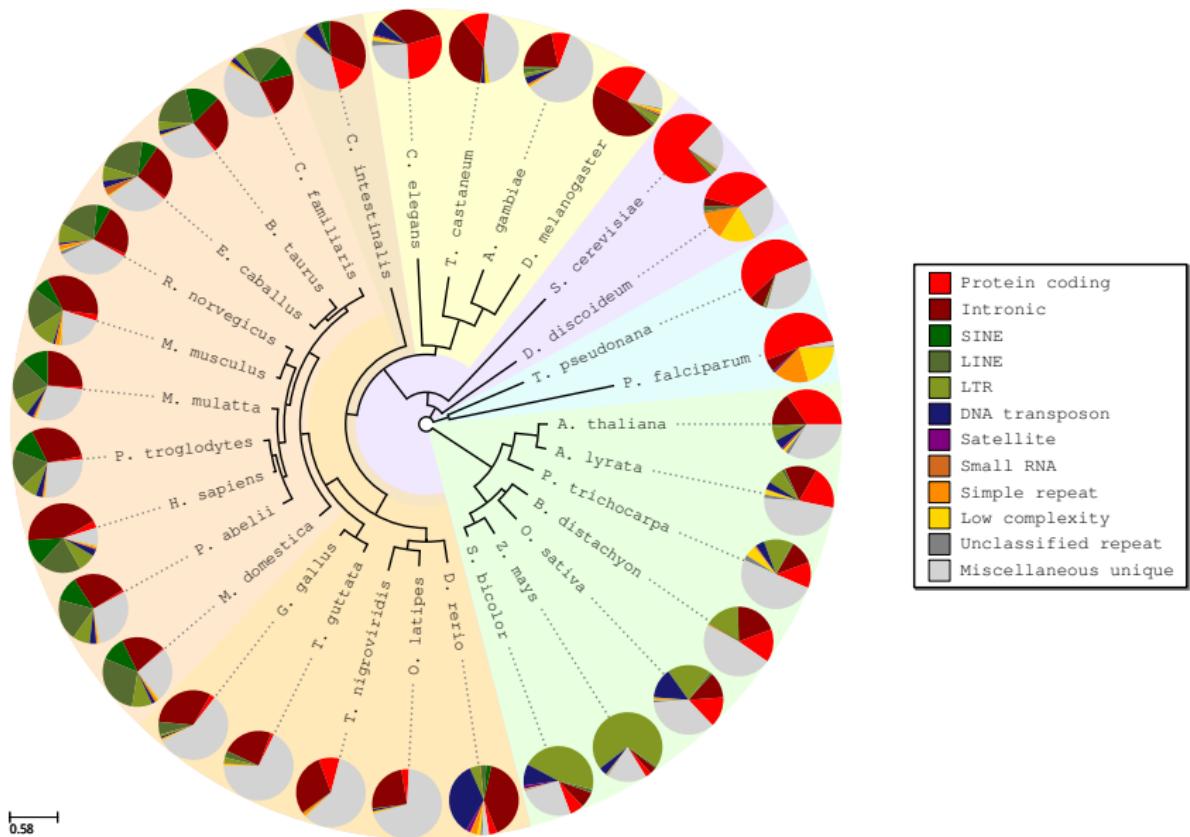
## 3 Ecology of genetic elements

- Background
- Results

## 4 Genomic study of selective pressures in group of genes

- Background
- Methodology
- Results

# Distribution of genetic species within biological species



# Non-random distribution of genetic species

- Test for the simplest hypotheses: **each of the genetic species (GSs) are found in equal proportions in each chromosome of a given genome**
- We generated 1,000 random genomes for each of the 31 biological species
- Test the differences between the GSs in the 1,000 shuffling and the observed value in 548 chromosomes
- In less than **4% of the cases GS were found in the expected proportion**

# Non-random distribution of genetic species

- Test for the simplest hypotheses: **each of the genetic species (GSs) are found in equal proportions in each chromosome of a given genome**
- We generated 1,000 random genomes for each of the 31 biological species
- Test the differences between the GSs in the 1,000 shuffling and the observed value in 548 chromosomes
- In less than **4% of the cases GS were found in the expected proportion**

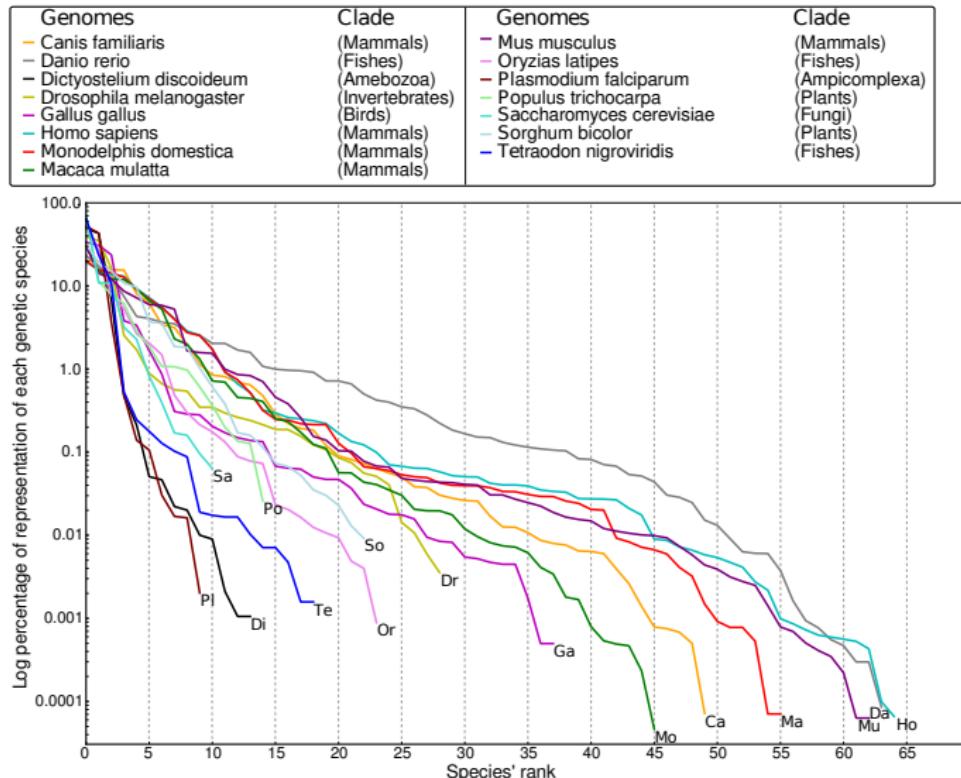
# Non-random distribution of genetic species

- Test for the simplest hypotheses: **each of the genetic species (GSs) are found in equal proportions in each chromosome of a given genome**
- We generated 1,000 random genomes for each of the 31 biological species
- Test the differences between the GSs in the 1,000 shuffling and the observed value in 548 chromosomes
- In less than **4% of the cases GS were found in the expected proportion**

# Non-random distribution of genetic species

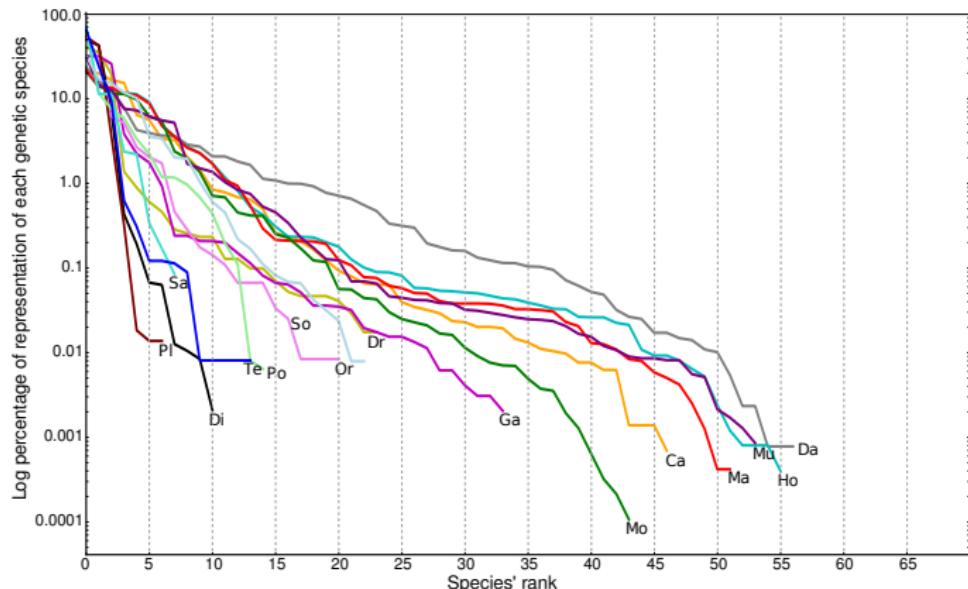
- Test for the simplest hypotheses: **each of the genetic species (GSs) are found in equal proportions in each chromosome of a given genome**
- We generated 1,000 random genomes for each of the 31 biological species
- Test the differences between the GSs in the 1,000 shuffling and the observed value in 548 chromosomes
- In less than **4% of the cases GS were found in the expected proportion**

# Counterbalanced GSs abundances in genomes

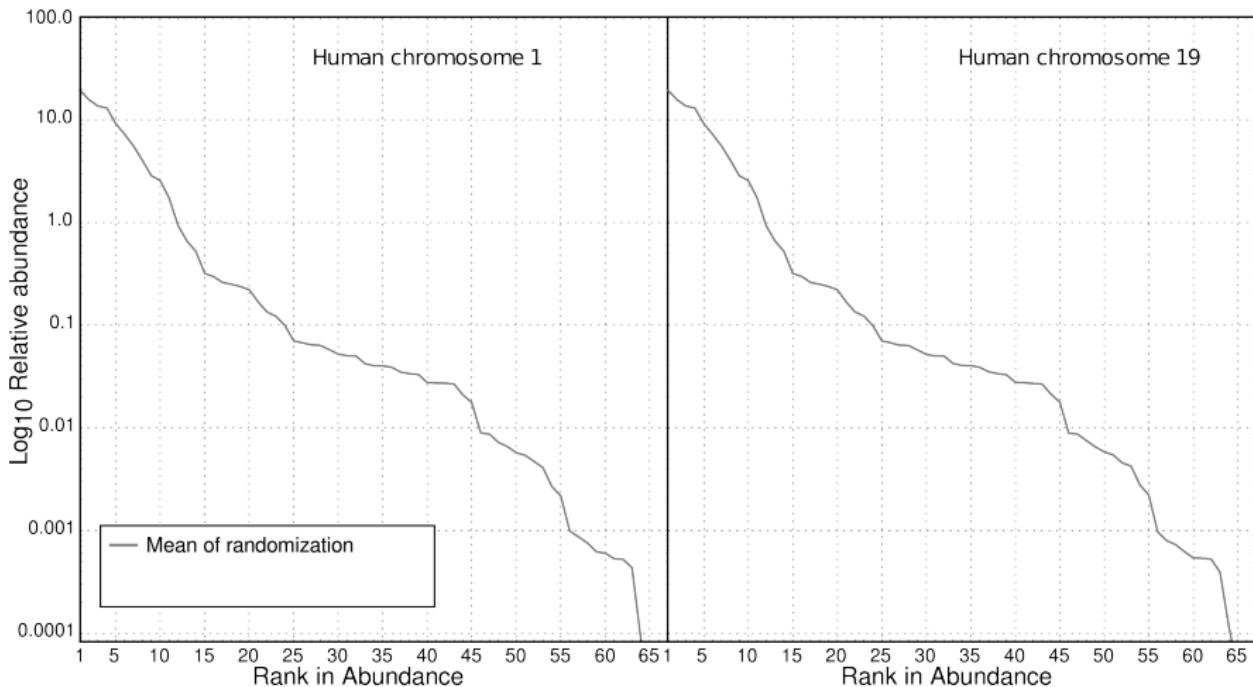


# Counterbalanced GSs abundances in Chromosomes

Genomes	Chrm.	Clade	Genomes	Chrm.	Clade
Canis familiaris	1	(Mammals)	Mus musculus	1	(Mammals)
Danio rerio	7	(Fishes)	Oryzias latipes	3	(Fishes)
Dictyostelium discoideum	2	(Amebozoa)	Plasmodium falciparum	14	(Apicomplexa)
Drosophila melanogaster	3R	(Invertebrates)	Populus trichocarpa	1	(Plants)
Gallus gallus	1	(Birds)	Saccharomyces cerevisiae	IV	(Fungi)
Homo sapiens	1	(Mammals)	Sorghum bicolor	2	(Plants)
Monodelphis domestica	1	(Mammals)	Tetraodon nigroviridis	1	(Fishes)
Macaca mulatta	1	(Mammals)			

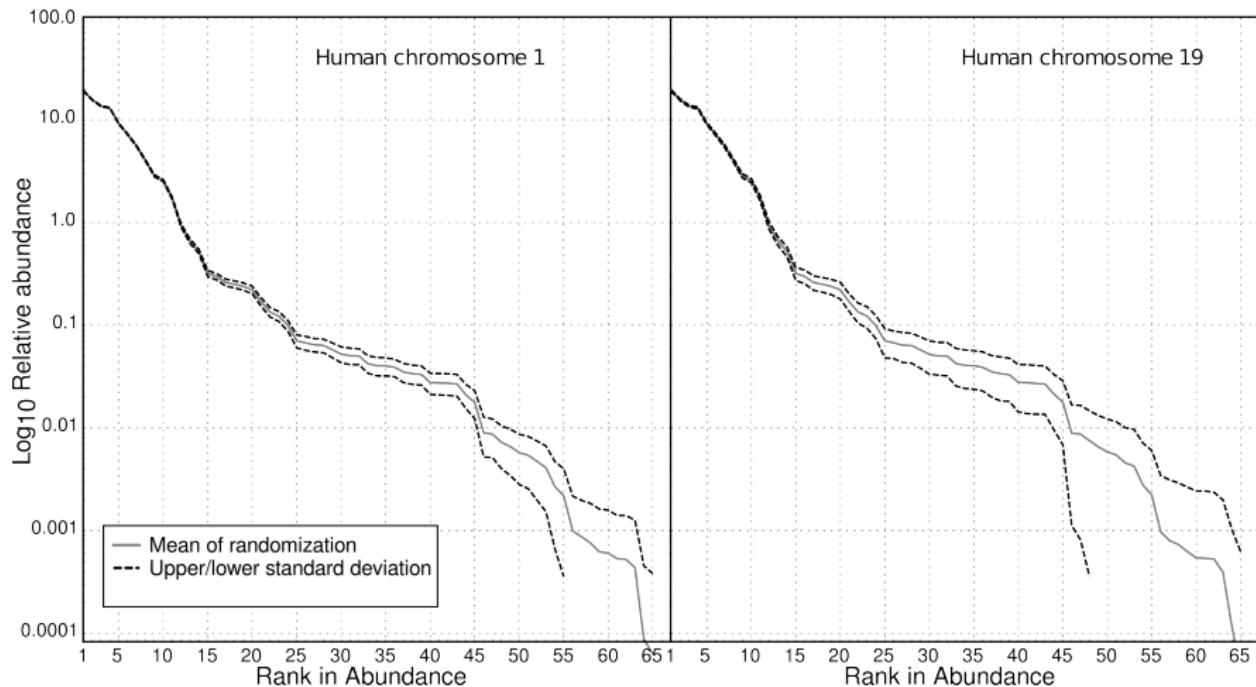


# Relative species abundance: observed vs expected



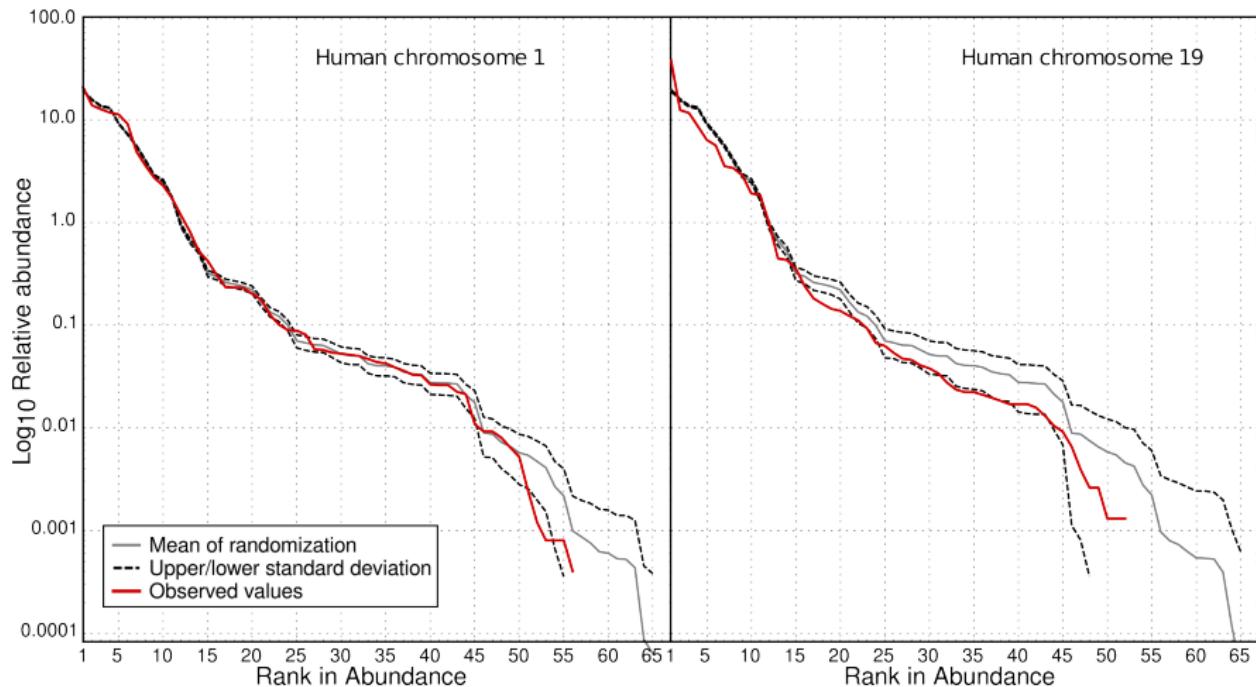
- Kolmogorov test red vs grey lines. 76 out of 548 chromosomes showing different distributions → **86% of chromosomes fitting**

# Relative species abundance: observed vs expected



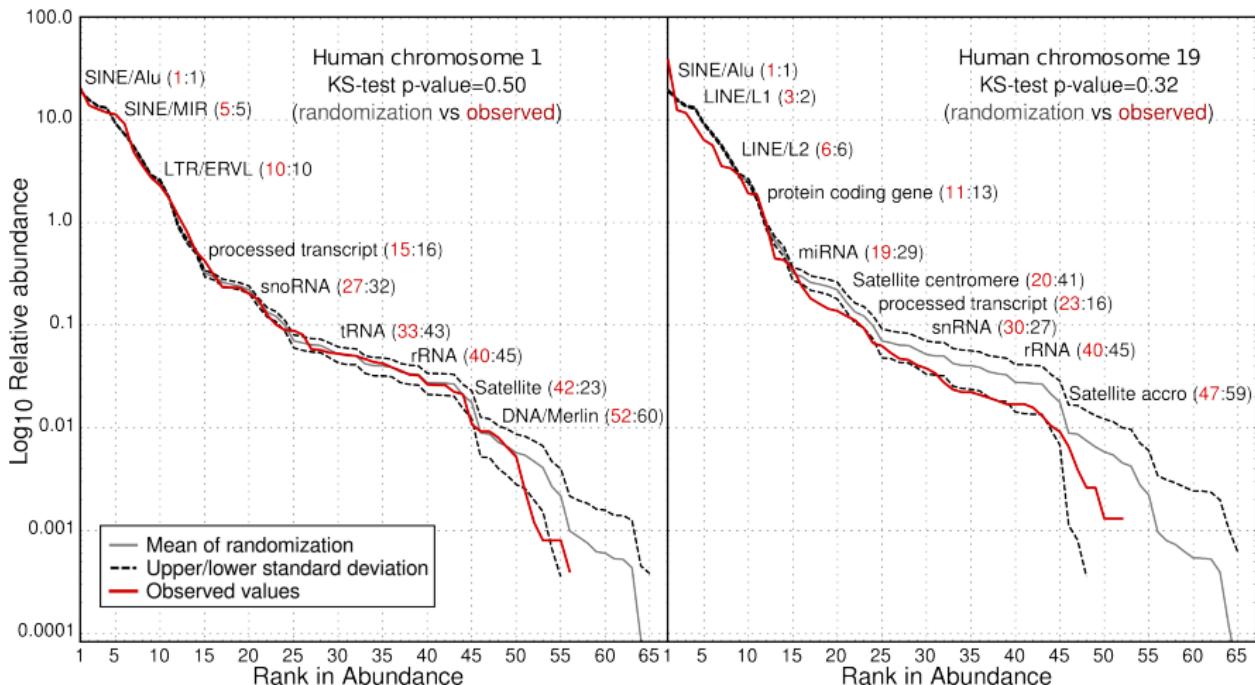
- Kolmogorov test red vs grey lines. 76 out of 548 chromosomes showing different distributions → **86% of chromosomes fitting**

# Relative species abundance: observed vs expected



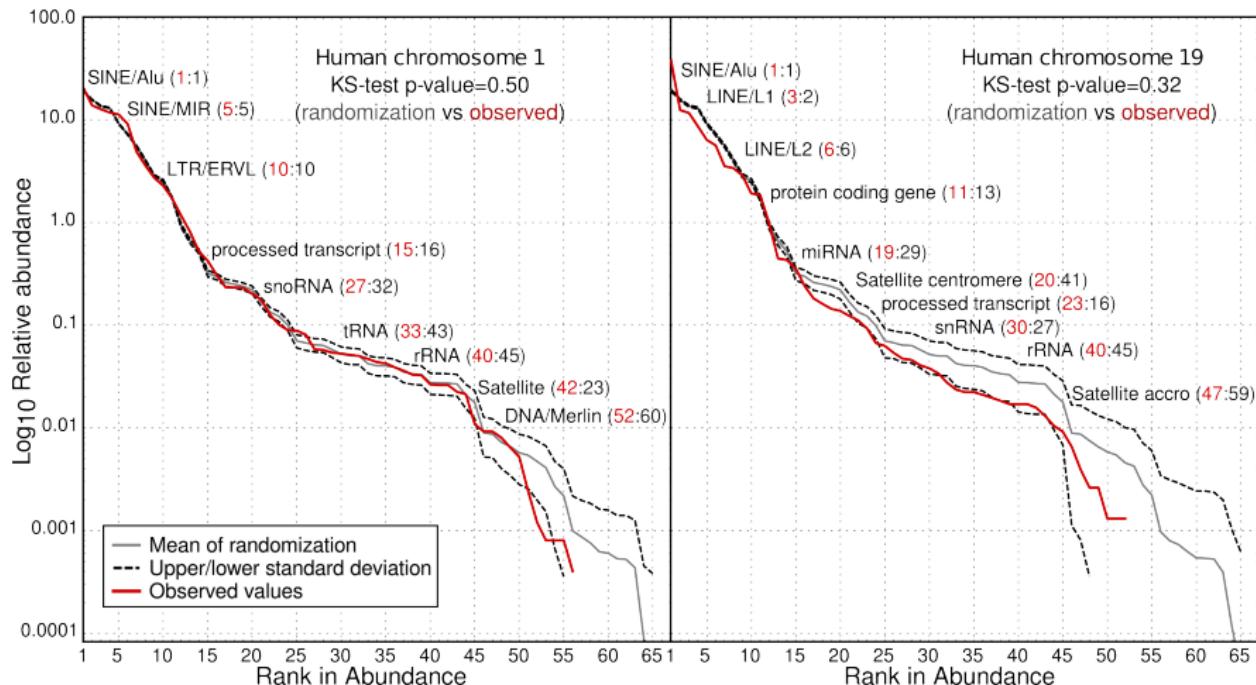
- Kolmogorov test red vs grey lines. 76 out of 548 chromosomes showing different distributions → **86% of chromosomes fitting**

# Relative species abundance: observed vs expected



- Kolmogorov test **red** vs grey lines. 76 out of 548 chromosomes showing different distributions → **86% of chromosomes fitting**

# Relative species abundance: observed vs expected



- Kolmogorov test **red** vs **grey** lines. 76 out of 548 chromosomes showing different distributions → **86% of chromosomes fitting**

# Neutrality of GSs distribution in chromosomes

Unified Neutral Theory of Biodiversity and Biogeography for genetic species

- We test the fit of GSs' distributions in the UNTB

Species	Chr	J	S	$\theta$	m	P-val
<i>T. castaneum</i>	7	7,865	18	2.12	—	<b>0.01</b>
<i>A. gambiae</i>	X	21,215	42	6.97	0.037	<b>0.03</b>
<i>G. gallus</i>	9	6,621	32	4.28	—	<b>0.05</b>
<i>D. melanogaster</i>	X	20,787	26	2.86	—	<b>0.09</b>
<i>T. nigroviridis</i>	3	8,505	17	1.96	—	<b>0.11</b>
<i>M. musculus</i>	14	143,018	59	6.58	0.149	<b>0.15</b>
<i>P. trichocarpa</i>	2	32,946	15	2.23	0.009	<b>0.16</b>
<i>O. latipes</i>	19	7,223	21	2.57	—	<b>0.17</b>
<i>H. sapiens</i>	17	93,105	52	6.51	0.065	<b>0.18</b>
<i>M. mulatta</i>	16	84,626	50	5.95	0.119	<b>0.23</b>
<i>S. cerevisiae</i>	II	640	9	1.37	—	<b>0.26</b>
<i>D. discoideum</i>	1	26,650	14	1.36	—	<b>0.27</b>
<i>D. rerio</i>	1	105,305	56	8.17	0.016	<b>0.29</b>
<i>C. familiaris</i>	1	144,103	47	5.28	0.093	<b>0.32</b>
<i>P. falciparum</i>	13	18,738	10	0.95	—	<b>0.38</b>
<i>M. domestica</i>	2	675,788	44	4.46	0.031	<b>0.43</b>
<i>S. bicolor</i>	1	37,626	23	2.86	0.067	<b>0.68</b>

- Neutrality can not be rejected for any of the chromosome studied

# Neutrality of GSs distribution in chromosomes

Unified Neutral Theory of Biodiversity and Biogeography for genetic species

- We test the fit of GSs' distributions in the UNTB

Species	Chr	J	S	$\theta$	m	P-val
<i>T. castaneum</i>	7	7,865	18	2.12	—	0.01
<i>A. gambiae</i>	X	21,215	42	6.97	0.037	0.03
<i>G. gallus</i>	9	6,621	32	4.28	—	0.05
<i>D. melanogaster</i>	X	20,787	26	2.86	—	0.09
<i>T. nigroviridis</i>	3	8,505	17	1.96	—	0.11
<i>M. musculus</i>	14	143,018	59	6.58	0.149	0.15
<i>P. trichocarpa</i>	2	32,946	15	2.23	0.009	0.16
<i>O. latipes</i>	19	7,223	21	2.57	—	0.17
<i>H. sapiens</i>	17	93,105	52	6.51	0.065	0.18
<i>M. mulatta</i>	16	84,626	50	5.95	0.119	0.23
<i>S. cerevisiae</i>	II	640	9	1.37	—	0.26
<i>D. discoideum</i>	1	26,650	14	1.36	—	0.27
<i>D. rerio</i>	1	105,305	56	8.17	0.016	0.29
<i>C. familiaris</i>	1	144,103	47	5.28	0.093	0.32
<i>P. falciparum</i>	13	18,738	10	0.95	—	0.38
<i>M. domestica</i>	2	675,788	44	4.46	0.031	0.43
<i>S. bicolor</i>	1	37,626	23	2.86	0.067	0.68

- Neutrality can not be rejected for any of the chromosome studied

# Neutrality of GSs distribution in chromosomes

Unified Neutral Theory of Biodiversity and Biogeography for genetic species

- We test the fit of GSs' distributions in the UNTB

Species	Chr	J	S	$\theta$	m	P-val
<i>T. castaneum</i>	7	7,865	18	2.12	—	0.01
<i>A. gambiae</i>	X	21,215	42	6.97	0.037	0.03
<i>G. gallus</i>	9	6,621	32	4.28	—	0.05
<i>D. melanogaster</i>	X	20,787	26	2.86	—	0.09
<i>T. nigroviridis</i>	3	8,505	17	1.96	—	0.11
<i>M. musculus</i>	14	143,018	59	6.58	0.149	0.15
<i>P. trichocarpa</i>	2	32,946	15	2.23	0.009	0.16
<i>O. latipes</i>	19	7,223	21	2.57	—	0.17
<i>H. sapiens</i>	17	93,105	52	6.51	0.065	0.18
<i>M. mulatta</i>	16	84,626	50	5.95	0.119	0.23
<i>S. cerevisiae</i>	II	640	9	1.37	—	0.26
<i>D. discoideum</i>	1	26,650	14	1.36	—	0.27
<i>D. rerio</i>	1	105,305	56	8.17	0.016	0.29
<i>C. familiaris</i>	1	144,103	47	5.28	0.093	0.32
<i>P. falciparum</i>	13	18,738	10	0.95	—	0.38
<i>M. domestica</i>	2	675,788	44	4.46	0.031	0.43
<i>S. bicolor</i>	1	37,626	23	2.86	0.067	0.68

- Neutrality can not be rejected for any of the chromosome studied

# Neutrality of GSs distribution in chromosomes

Unified Neutral Theory of Biodiversity and Biogeography for genetic species

- We test the fit of GSs' distributions in the UNTB

Species	Chr	J	S	$\theta$	m	P-val
<i>T. castaneum</i>	7	7,865	18	2.12	—	0.01
<i>A. gambiae</i>	X	21,215	42	6.97	0.037	0.03
<i>G. gallus</i>	9	6,621	32	4.28	—	0.05
<i>D. melanogaster</i>	X	20,787	26	2.86	—	0.09
<i>T. nigroviridis</i>	3	8,505	17	1.96	—	0.11
<i>M. musculus</i>	14	143,018	59	6.58	0.149	0.15
<i>P. trichocarpa</i>	2	32,946	15	2.23	0.009	0.16
<i>O. latipes</i>	19	7,223	21	2.57	—	0.17
<i>H. sapiens</i>	17	93,105	52	6.51	0.065	0.18
<i>M. mulatta</i>	16	84,626	50	5.95	0.119	0.23
<i>S. cerevisiae</i>	II	640	9	1.37	—	0.26
<i>D. discoideum</i>	1	26,650	14	1.36	—	0.27
<i>D. rerio</i>	1	105,305	56	8.17	0.016	0.29
<i>C. familiaris</i>	1	144,103	47	5.28	0.093	0.32
<i>P. falciparum</i>	13	18,738	10	0.95	—	0.38
<i>M. domestica</i>	2	675,788	44	4.46	0.031	0.43
<i>S. bicolor</i>	1	37,626	23	2.86	0.067	0.68

- Neutrality can not be rejected for any of the chromosome studied

# Neutrality of GSs distribution in chromosomes

Unified Neutral Theory of Biodiversity and Biogeography for genetic species

- We test the fit of GSs' distributions in the UNTB

Species	Chr	J	S	$\theta$	m	P-val
<i>T. castaneum</i>	7	7,865	18	2.12	—	<b>0.01</b>
<i>A. gambiae</i>	X	21,215	42	6.97	0.037	<b>0.03</b>
<i>G. gallus</i>	9	6,621	32	4.28	—	<b>0.05</b>
<i>D. melanogaster</i>	X	20,787	26	2.86	—	<b>0.09</b>
<i>T. nigroviridis</i>	3	8,505	17	1.96	—	<b>0.11</b>
<i>M. musculus</i>	14	143,018	59	6.58	0.149	<b>0.15</b>
<i>P. trichocarpa</i>	2	32,946	15	2.23	0.009	<b>0.16</b>
<i>O. latipes</i>	19	7,223	21	2.57	—	<b>0.17</b>
<i>H. sapiens</i>	17	93,105	52	6.51	0.065	<b>0.18</b>
<i>M. mulatta</i>	16	84,626	50	5.95	0.119	<b>0.23</b>
<i>S. cerevisiae</i>	II	640	9	1.37	—	<b>0.26</b>
<i>D. discoideum</i>	1	26,650	14	1.36	—	<b>0.27</b>
<i>D. rerio</i>	1	105,305	56	8.17	0.016	<b>0.29</b>
<i>C. familiaris</i>	1	144,103	47	5.28	0.093	<b>0.32</b>
<i>P. falciparum</i>	13	18,738	10	0.95	—	<b>0.38</b>
<i>M. domestica</i>	2	675,788	44	4.46	0.031	<b>0.43</b>
<i>S. bicolor</i>	1	37,626	23	2.86	0.067	<b>0.68</b>

- Neutrality can not be rejected for any of the chromosome studied

# Neutrality of GSs distribution in chromosomes

Unified Neutral Theory of Biodiversity and Biogeography for genetic species

- We test the fit of GSs' distributions in the UNTB

Species	Chr	J	S	$\theta$	m	P-val
<i>T. castaneum</i>	7	7,865	18	2.12	—	<b>0.01</b>
<i>A. gambiae</i>	X	21,215	42	6.97	0.037	<b>0.03</b>
<i>G. gallus</i>	9	6,621	32	4.28	—	<b>0.05</b>
<i>D. melanogaster</i>	X	20,787	26	2.86	—	<b>0.09</b>
<i>T. nigroviridis</i>	3	8,505	17	1.96	—	<b>0.11</b>
<i>M. musculus</i>	14	143,018	59	6.58	0.149	<b>0.15</b>
<i>P. trichocarpa</i>	2	32,946	15	2.23	0.009	<b>0.16</b>
<i>O. latipes</i>	19	7,223	21	2.57	—	<b>0.17</b>
<i>H. sapiens</i>	17	93,105	52	6.51	0.065	<b>0.18</b>
<i>M. mulatta</i>	16	84,626	50	5.95	0.119	<b>0.23</b>
<i>S. cerevisiae</i>	II	640	9	1.37	—	<b>0.26</b>
<i>D. discoideum</i>	1	26,650	14	1.36	—	<b>0.27</b>
<i>D. rerio</i>	1	105,305	56	8.17	0.016	<b>0.29</b>
<i>C. familiaris</i>	1	144,103	47	5.28	0.093	<b>0.32</b>
<i>P. falciparum</i>	13	18,738	10	0.95	—	<b>0.38</b>
<i>M. domestica</i>	2	675,788	44	4.46	0.031	<b>0.43</b>
<i>S. bicolor</i>	1	37,626	23	2.86	0.067	<b>0.68</b>

- Neutrality can not be rejected for any of the chromosome studied

# Index

## 1 Introduction

- Evolution only makes sense in light of neutrality
- Overview

## 2 Random-like structure of DNA

- Background
- Methodology
- Results

## 3 Ecology of genetic elements

- Background
- Results

## 4 Genomic study of selective pressures in group of genes

- **Background**
- Methodology
- Results

## The single gene approach

- Functional categories enriched in **PSGs** ( $\omega = \frac{dN}{dS} > 1$ )
- No set of PSGs were found to be significantly enriched

## The single gene approach

- Functional categories enriched in **PSGs** ( $\omega = \frac{dN}{dS} > 1$ )
- No set of PSGs were found to be significantly enriched

# The single gene approach

- Functional categories enriched in **PSGs** ( $\omega = \frac{dN}{dS} > 1$ )
- No set of PSGs were found to be significantly enriched



[Arbiza et al. - PLoS Computational Biology (2006), Clark et al. - Science (2003), Nielsen et al. - PLoS Biology (2005)]

# Index

## 1 Introduction

- Evolution only makes sense in light of neutrality
- Overview

## 2 Random-like structure of DNA

- Background
- Methodology
- Results

## 3 Ecology of genetic elements

- Background
- Results

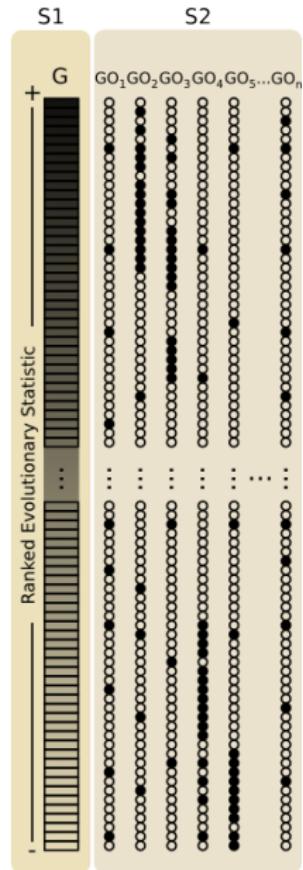
## 4 Genomic study of selective pressures in group of genes

- Background
- Methodology**
- Results

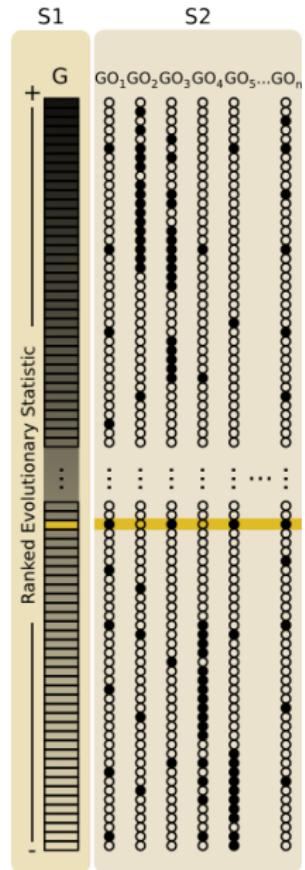
# Gene-Set Selection Analysis (GSSA)



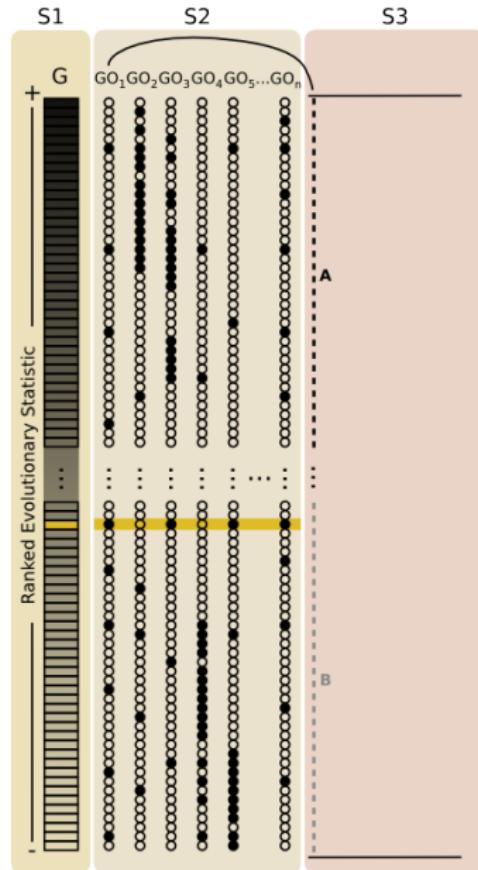
# Gene-Set Selection Analysis (GSSA)



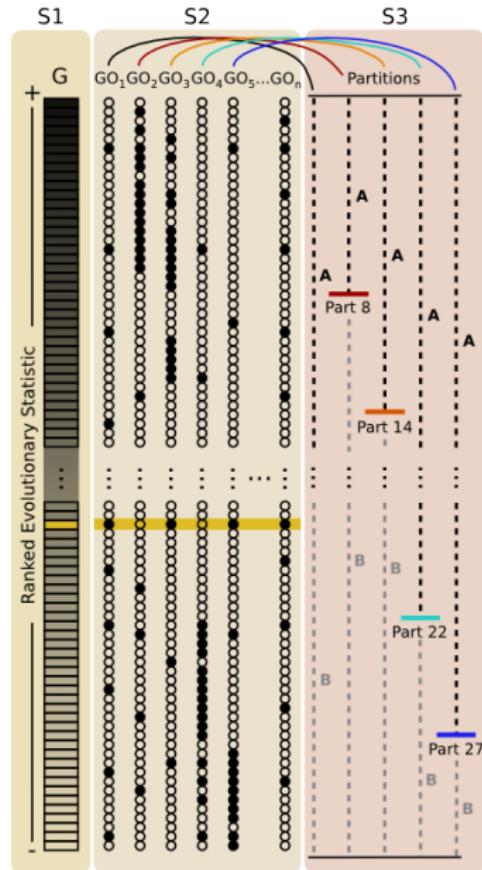
# Gene-Set Selection Analysis (GSSA)



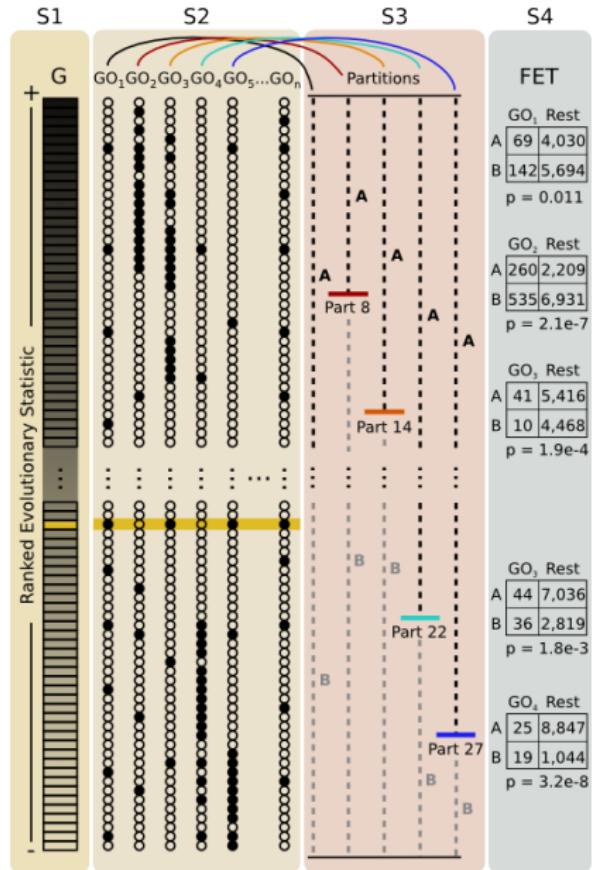
# Gene-Set Selection Analysis (GSSA)



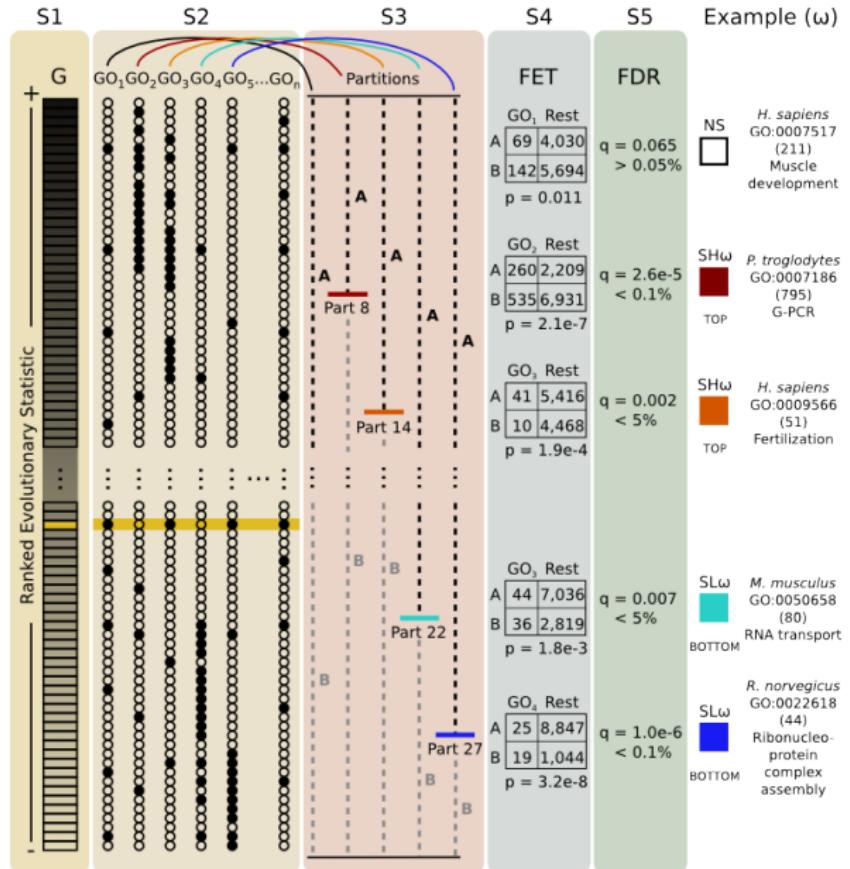
# Gene-Set Selection Analysis (GSSA)



# Gene-Set Selection Analysis (GSSA)



# Gene-Set Selection Analysis (GSSA)



# Index

## 1 Introduction

- Evolution only makes sense in light of neutrality
- Overview

## 2 Random-like structure of DNA

- Background
- Methodology
- Results

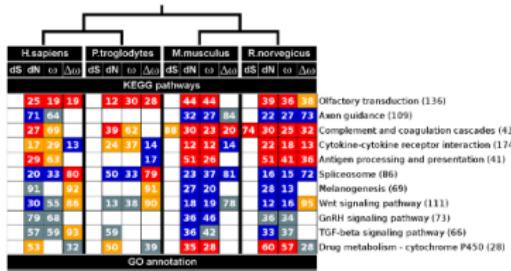
## 3 Ecology of genetic elements

- Background
- Results

## 4 Genomic study of selective pressures in group of genes

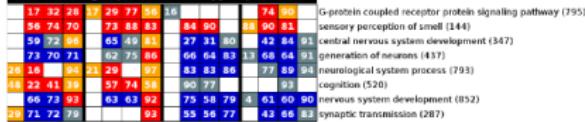
- Background
- Methodology
- Results

# GSSA – Summary results



## GO annotation

## Neurological process and sensory perception



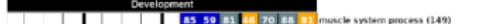
## Immunity and defense response



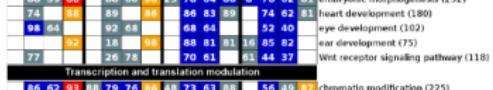
## Reproduction



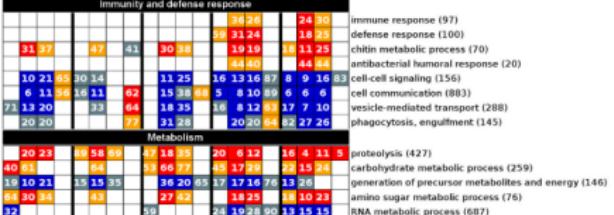
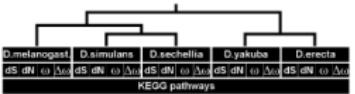
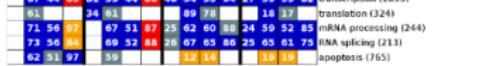
## Metabolism



## Development



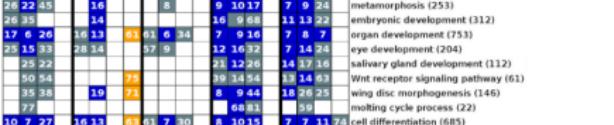
## Transcription and translation modulation



## Metabolism



## Reproduction



## Transcription and translation modulation



## RNA processing (187)

## transcription (541)

## translation (264)

## gene expression (1000)

## metamorphosis (253)

## embryonic development (312)

## organ development (753)

## eye development (204)

## salivary gland development (112)

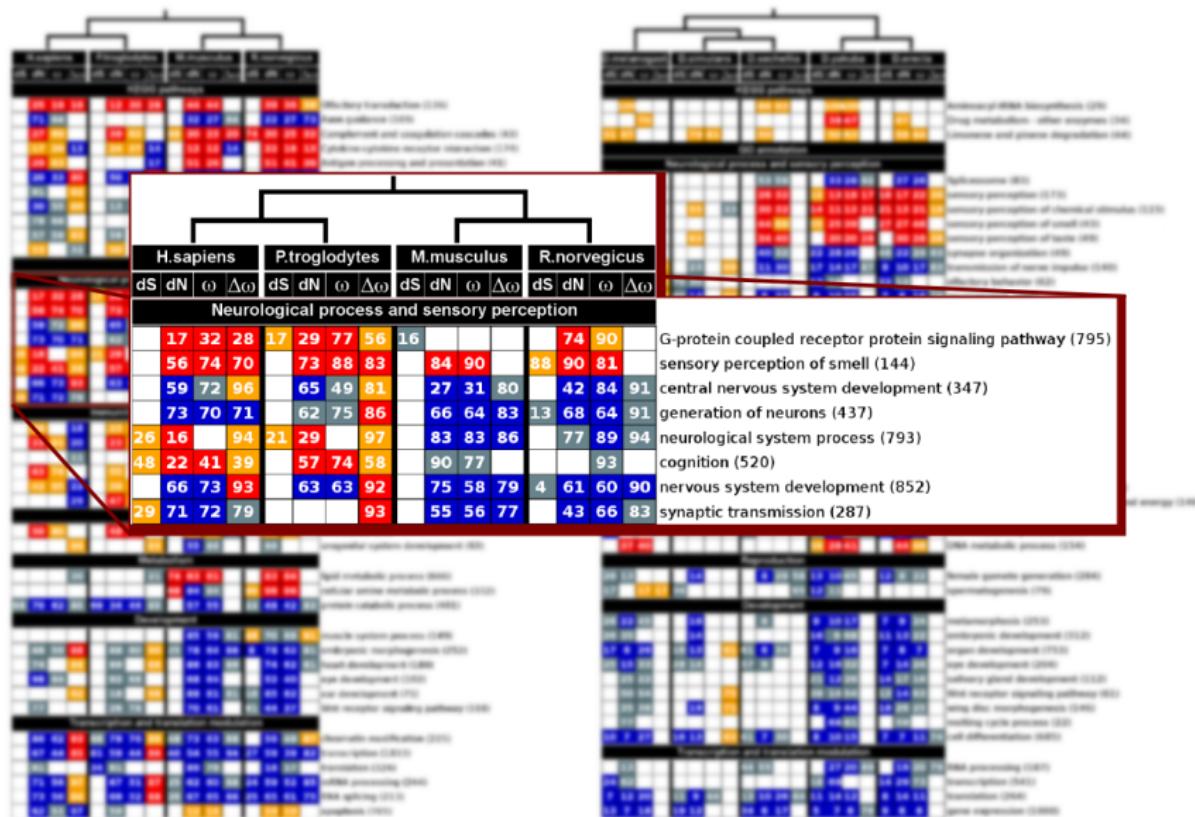
## Wnt receptor signaling pathway (61)

## wing disc morphogenesis (146)

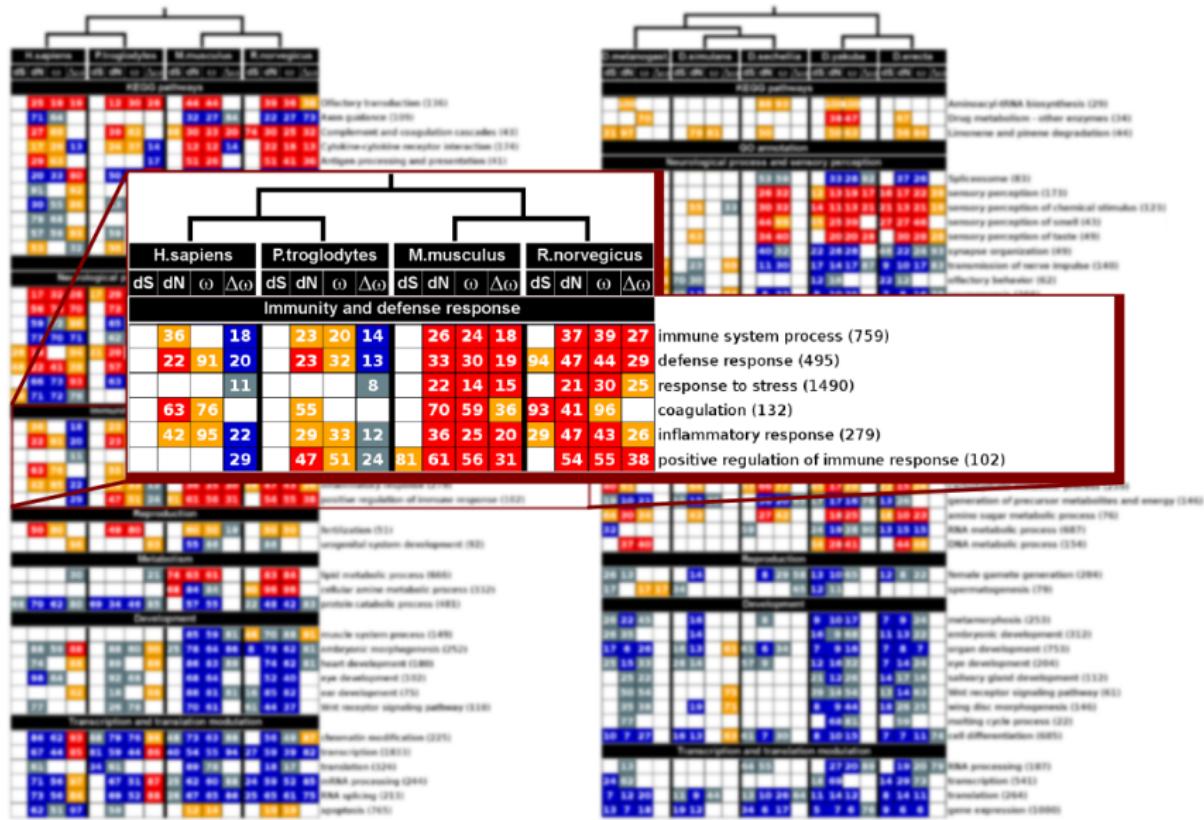
## molting cycle process (22)

## cell differentiation (685)

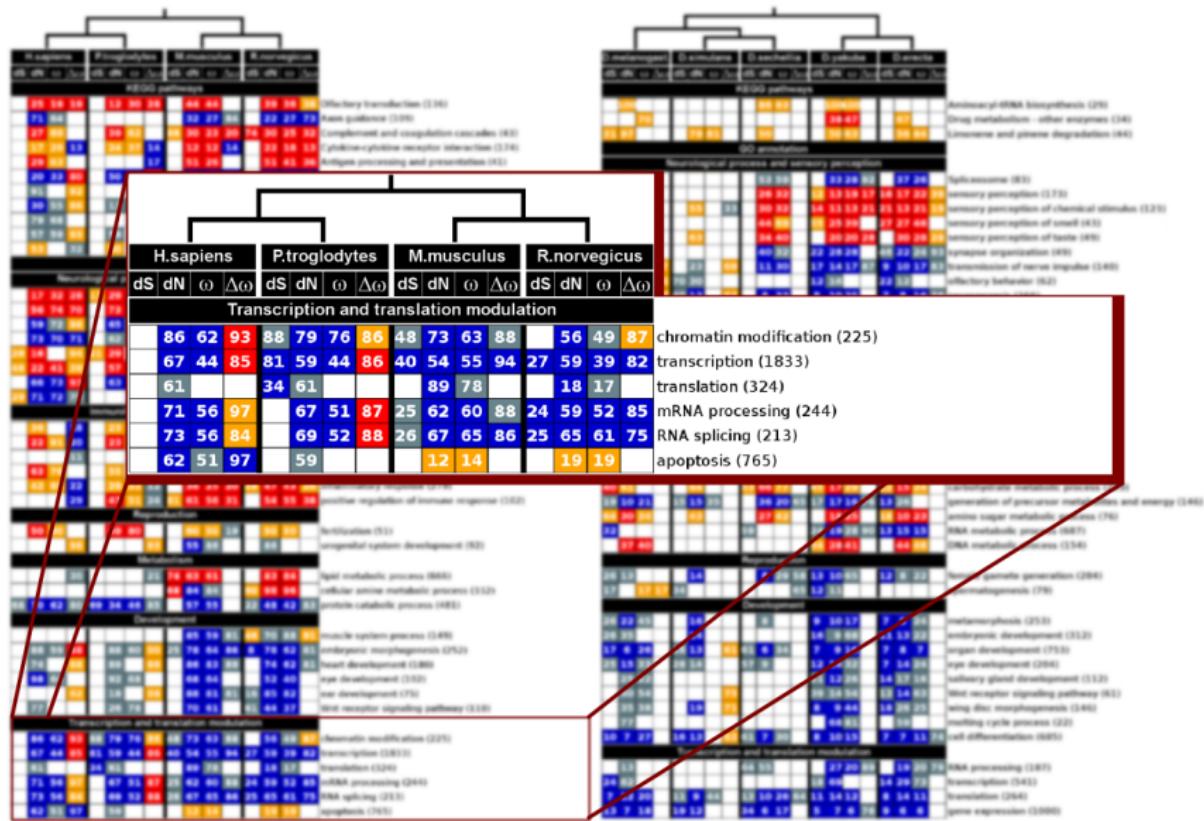
## GSSA – Summary results



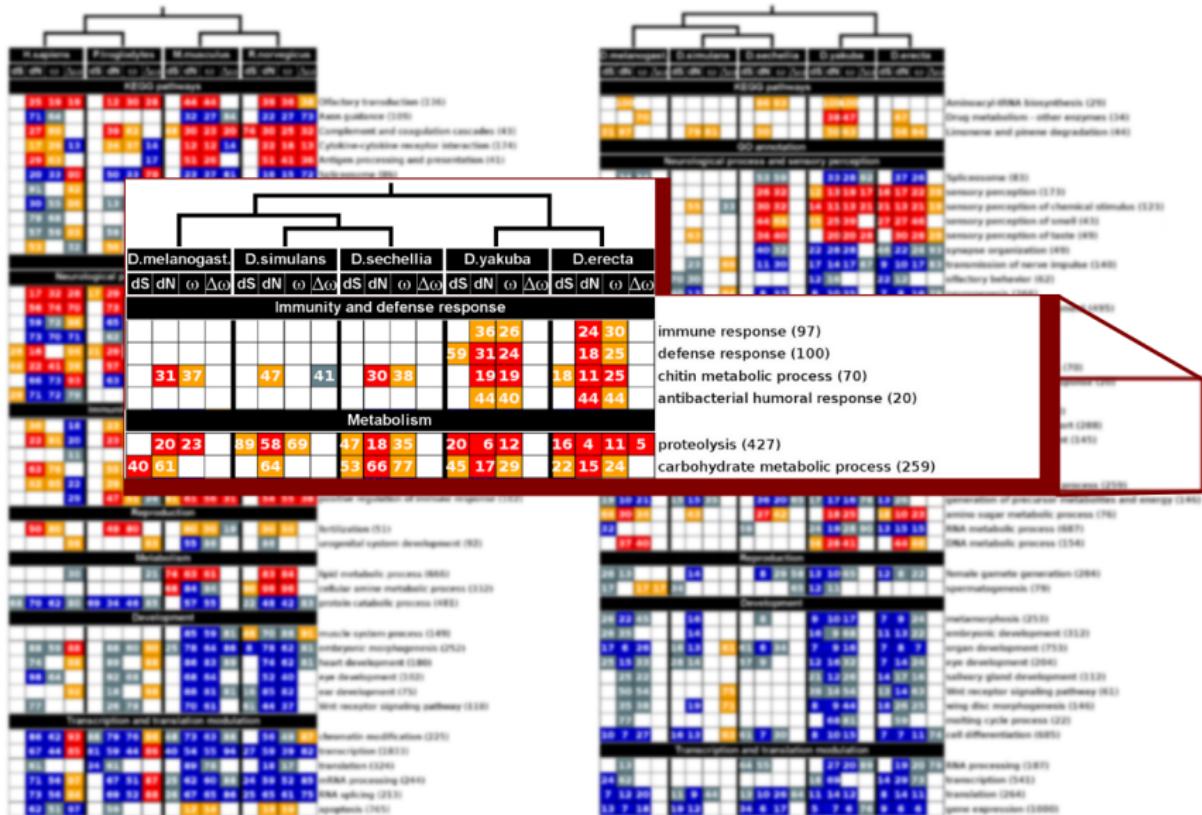
# GSSA – Summary results



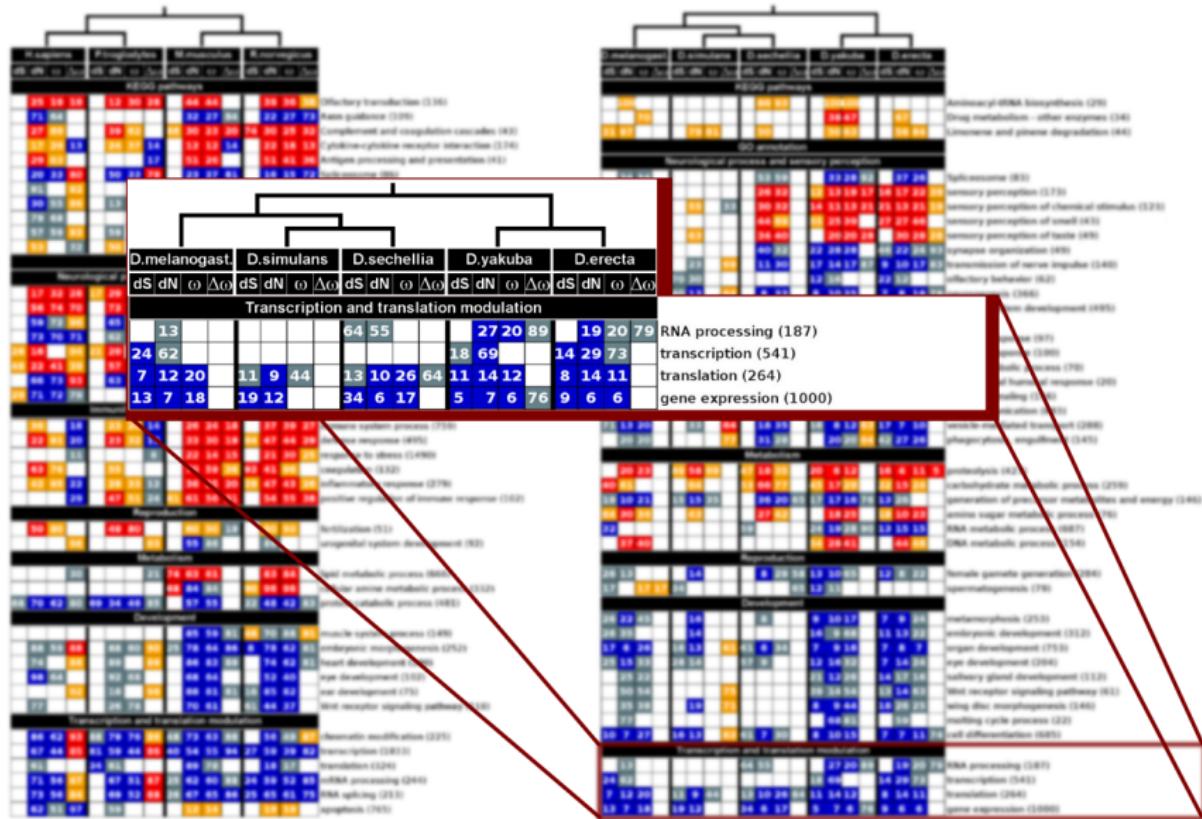
## GSSA – Summary results



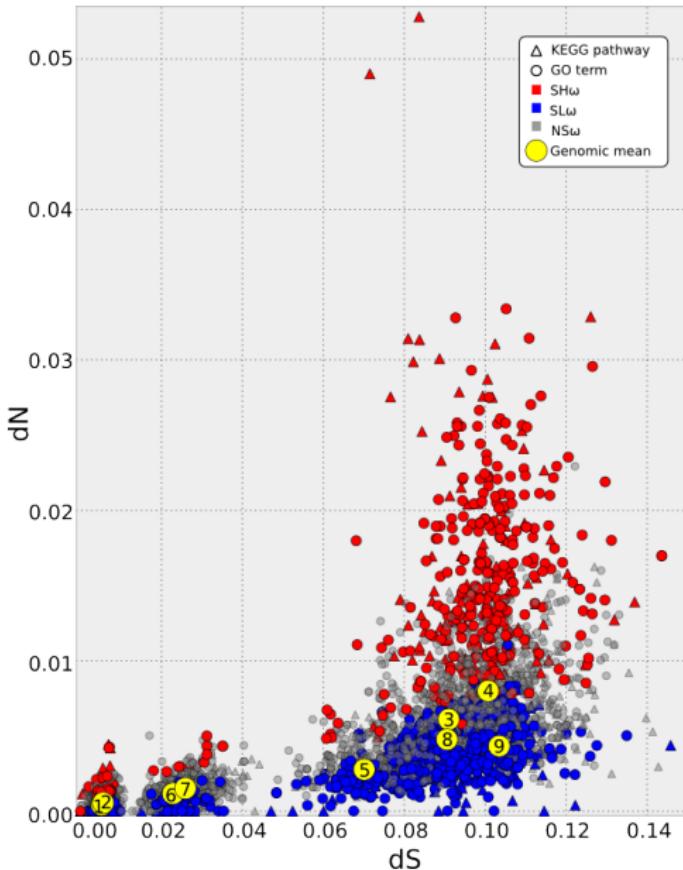
# GSSA – Summary results



# GSSA – Summary results

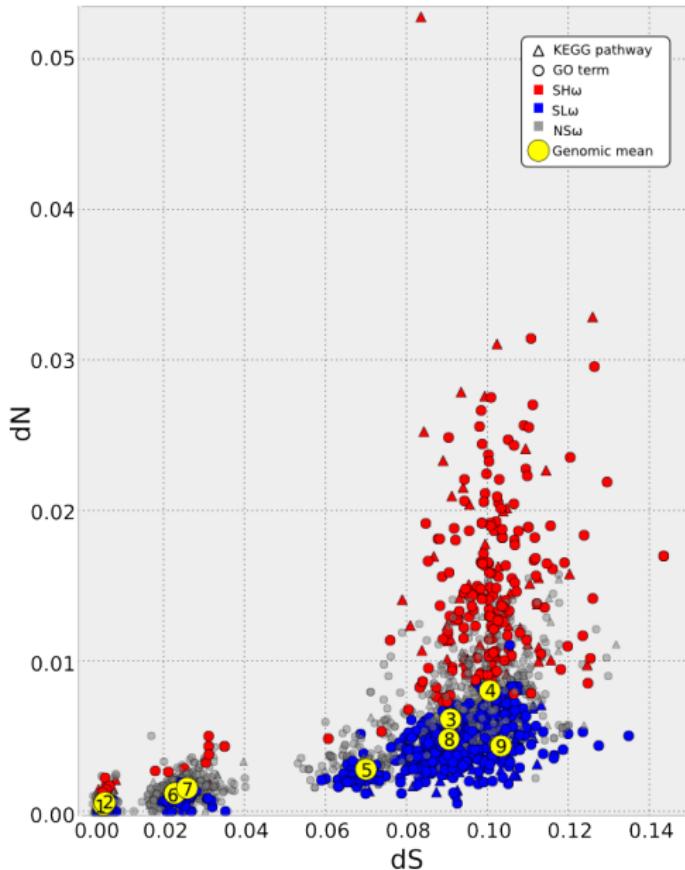


# PSGs within GSSA results



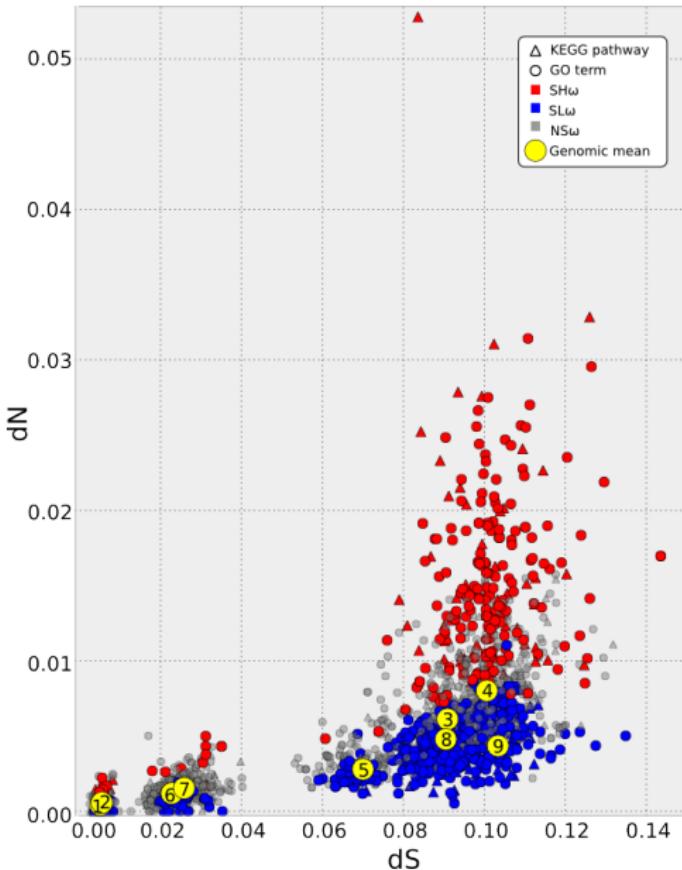
- PSGs are largely present in categories with  $SH\omega$ ,  $SL\omega$  and  $NS\omega$ .
- PSGs were enriched in  $SH\omega$  categories
- However not significant for primates

# PSGs within GSSA results



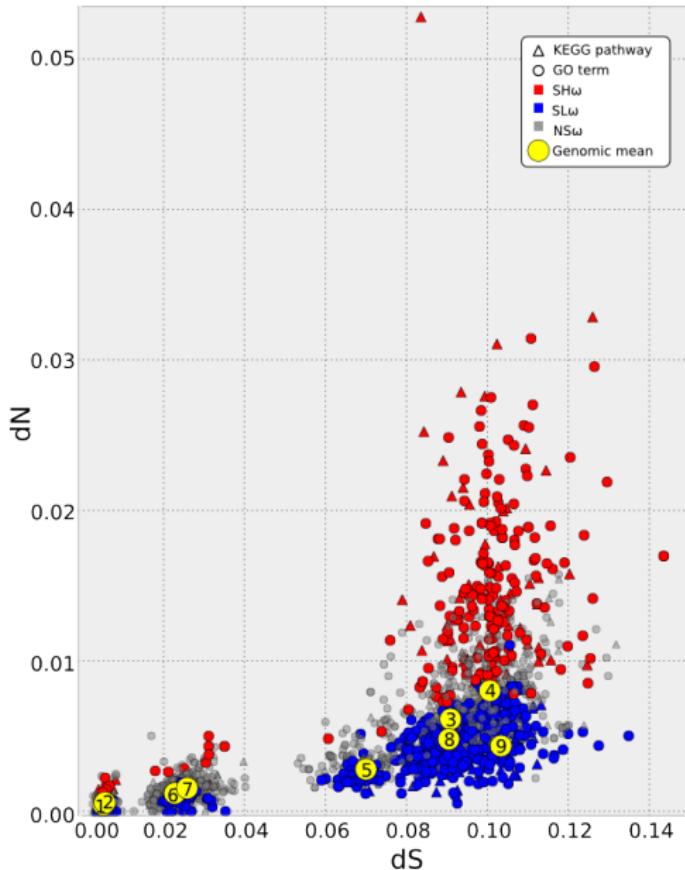
- PSGs are largely present in categories with  $SH\omega$ ,  $SL\omega$  and  $NS\omega$ .
- PSGs were enriched in  $SH\omega$  categories
- However not significant for primates

# PSGs within GSSA results



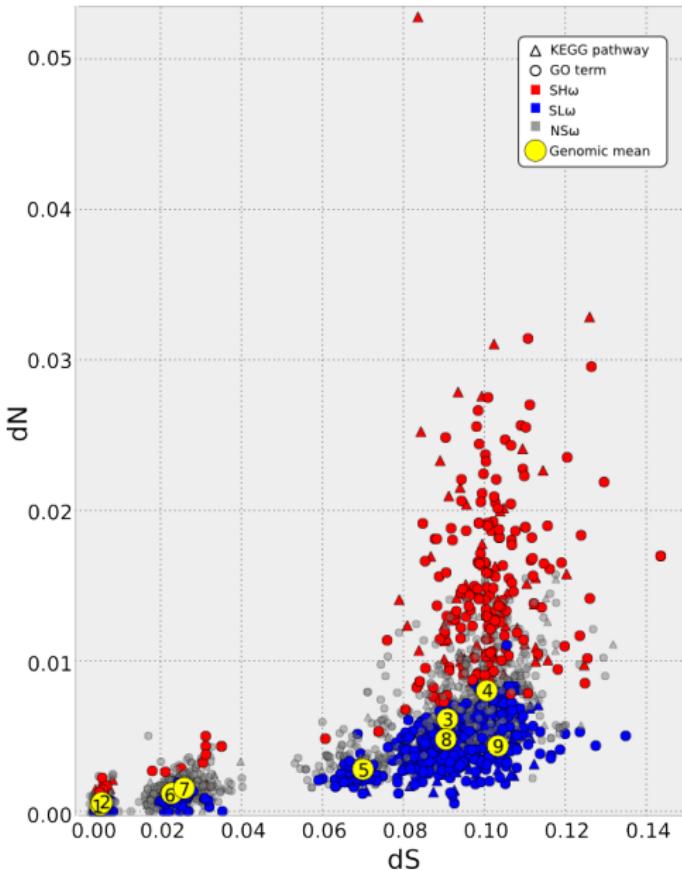
- PSGs are largely present in categories with  $SH\omega$ ,  $SL\omega$  and  $NS\omega$ .
- PSGs were enriched in  $SH\omega$  categories
- However not significant for primates

# PSGs within GSSA results



- PSGs are largely present in categories with  $SH\omega$ ,  $SL\omega$  and  $NS\omega$ .
- PSGs were enriched in  $SH\omega$  categories
- However not significant for primates

# PSGs within GSSA results



- PSGs are largely present in categories with  $SH\omega$ ,  $SL\omega$  and  $NS\omega$ .
- PSGs were enriched in  $SH\omega$  categories
- However not significant for primates

# Conclusions

- 1 In the whole diversity of life, from viruses to mammals, informational content of the genomes exhibits quasi-maximum values. Only dramatic changes such as polyploidization events, or strong biases in nucleotide contents, are able to lower genome entropy.
- 2 According to the observed universal adjustment of genomes to maximum complexity, we hypothesize that increases in biological complexity are the consequence of genome expansions events through duplications or polyploidization.
- 3 Similarly to the biological species in ecosystems, eukaryotic genomes present an heterogeneous distribution of families or “species” of genetic elements: a few are very abundant, others quite frequent, and the majority rare.
- 4 Likewise for ecological species-area correlation graph, we observe that, along a great diversity of eukaryote genomes, the total number of genetic species in chromosomes is proportional to chromosome length.
- 5 The distributions and abundances of families of genetic elements in eukaryotic genomes, either functional or repetitive, follow the expectations of the unified neutral theory of biodiversity (UNTB).

# Conclusions

- 6 GSSA allows testing for functional biases within fast or slowly-evolving genes. This methodology successfully identified all previously reported candidate functional categories as important targets of natural selection. Moreover, given that the GSSA is not limited by the compulsory presence of positively-selected genes, it extended the list of phenotypical targets to previously undetectable ones.
- 7 Genes under positive selection were found to be present in functional categories evolving rapidly, slowly or without a significant trend ( $\omega$ ). However a significant bias towards fast evolving categories was found in rodents and *Drosophila*. Regarding to the even distribution found in primates, we hypothesize that it may be the result of small population sizes limiting the influence of natural selection; just as suggested by the theory of slightly deleterious mutations.
- 8 We believe that the role of genes under positive selection not only consists of the adaptive evolutionary changes of phenotypes but may also be related to other processes such as the adjustment to the deleterious mutations of other genes in a given network.
- 9 Throughout this thesis three primary bioinformatics tools were implemented with the idea of facilitating and extending future researches of the scientific community. These software are in line with the fields of ecology (of genomes), phylogeny, phylogenomics and the testing of evolutionary hypotheses.

Thank you :)

# Bibliography I

## thebibliography

-  LEONARDO ARBIZA, JOAQUÍN DOPAZO, AND HERNÁN DOPAZO, Positive selection, relaxation, and acceleration in the evolution of the human and chimp genome.  
*PLoS Computational Biology* **2**(4) (2006), p. e38.  
77
-  ANDREW G CLARK, STEPHEN GLANOWSKI, RASMUS NIELSEN, PAUL D THOMAS, ANISH KEJARIWAL, MELISSA A TODD, DAVID M TANENBAUM, DANIEL CIVELLO, FU LU, BRIAN MURPHY, STEVE FERRIERA, GARY WANG, XIANQUN ZHENG, THOMAS J WHITE, JOHN J SNINSKY, ET AL., Inferring nonneutral evolution from human-chimp-mouse orthologous gene trios.  
*Science* **302**(5652) (2003), 1960–3.  
77

## Bibliography II



T RYAN GREGORY, Synergy between sequence and size in large-scale genomics.

*Nature Reviews. Genetics* **6**(9) (2005), 699–708.

27, 28



STEPHEN P HUBBELL, *The Unified Neutral Theory of Biodiversity and Biogeography*.

Princeton University Press, 2001.

15, 16



ERIC S LANDER, L M LINTON, BRUCE W BIRREN, C NUSBAUM, MICHAEL C ZODY, JENNIFER BALDWIN, K DEVON, K DEWAR, M DOYLE, W FITZHUGH, R FUNKE, D GAGE, K HARRIS, A HEAFORD, J HOWLAND, ET AL., Initial sequencing and analysis of the human genome.

*Nature* **409**(6822) (2001), 860–921.

52

# Bibliography III



D. R. MADDISON AND K.-S. SCHULZ, The Tree of Life Web Project,  
<http://tolweb.org>, 2007.  
24, 25, 26



ANNE E MAGURRAN, *Ecological diversity and its measurement*.  
Cambridge University Press, 1988.

14



RASMUS NIELSEN, CARLOS BUSTAMANTE, ANDREW G CLARK,  
STEPHEN GLANOWSKI, TIMOTHY B SACKTON, MELISSA J HUBISZ, ADI  
FLEDEL-ALON, DAVID M TANENBAUM, DANIEL CIVELLO, THOMAS J  
WHITE, JOHN J SNINSKY, MARK D ADAMS, AND MICHELE CARGILL, A  
scan for positively selected genes in the genomes of humans and  
chimpanzees.

*PLoS Biology* 3(6) (2005), p. e170.

77

# Bibliography IV



ELLEN J PRITHAM, Transposable elements and factors influencing their success in eukaryotes.

*The Journal of Heredity* **100**(5) (2009), 648–55.

53



JAMES ROSINDELL, STEPHEN P HUBBELL, AND RAMPAL S ETIENNE,  
The unified neutral theory of biodiversity and biogeography at age ten.

*Trends in Ecology & Evolution* **26**(7) (2011), 340–8.

17