0.1 Measuring DNA complexity

0.1.1 The complexity ratio and complexity value

	Datating saguenes	<i>I</i> .	/11	$GG:$ \mathbf{BWT}	-		Chan	1:04		MODE
#	Rotating sequence					~	Char.			MTF
0	AACCTTCGTAGCATGG	0		G		G	a	t	c	0
$\frac{1}{2}$	ACCTTCGTAGCATGG A CCTTCGTAGCATGG AA	1 5		A T		G	A	\mathbf{T}	c	$\frac{1}{2}$
3	CTTCGTAGCATGG AAC	э 7		C		A T	g		с С	3
4	TTCGTAGCATGG AACC	15		G		C	a t	g	G	3
5	TCGTAGCATGG AACCT	13		A		G		a t	A	3
6	CGTAGCATGG AACCT	6		T		A	С	c	T	3
7	GTAGCATGG AACCTTC	11	➾	C	₽	Т	g a	g	C	3
8	TAGCATGG AACCTTCG	12	,	G		C	t	a	G	3
9	AGCATGG AACCTTCGT	2		A		G	c	t	A	3
10	GCATGG AACCTTCGTA	9		T		A	g	c	T	3
11	CATGG AACCTTCGTA G	4		C		Т	a	g	Ĉ	3
12	ATGG AACCTTCGTAGC	3		G		C	t	a	G	3
13	TGG AACCTTCGTAGCA	14		T		G	c	Т	a	2
14	GG AACCTTCGTAGCA T	10		A		T	g	c	Α	3
15	G AACCTTCGTAGCAT G	8		С		A	t	g	\mathbf{C}	3
able	eq) = E(MTF(BWT(seq))) = 2 0.1: CR explained by each three tables summarizes	xam	ple	е.						
able Thes	0.1: CR explained by ester three tables summarizes ber from which we will final	xam the s	ple tep	e. s needed oute Sha	l to	obt on's	ain th	e fir y. 1	nal s	equence ne table
able Thes numl	0.1: CR explained by este three tables summarizes ber from which we will final eft corresponds to the BWT	xam the s lly co	ple tep mp	e. s needed oute Sha nal seque	l to nnc	obt on's e is r	ain the entrop	e fir oy. 1 l seq	nal s) Tl uen	equence ne table tially (fi
able Thes numl the le	0.1: CR explained by estate three tables summarizes ber from which we will final eft corresponds to the BWT acter moved to back) result	xam the s ly co . Ori ing ir	tep mp igin	e. s needed oute Sha hal seque	l to nno ence	obton's e is r	ain the entrop otated as ma	e fir by. 1 l seq ny a	nal s) Th luen s ch	equence ne table tially (fi aracters
able Thes numl the le chara the s	0.1: CR explained by enter three tables summarizes therefore which we will final eff corresponds to the BWT eacter moved to back) result sequence. The resulting seconds	xam the s lly co . Ori ing ir	tep mp igir n di	es needed oute Shar all seque afferent s	l to nnce ence etrin	obton's e is r	tain the entrope otated as ma	e fir by. 1 l seq ny a xico	nal s) Th uen s ch grap	equence ne table tially (fir aracters phic ord
Thes numl the le	0.1: CR explained by estate three tables summarizes ber from which we will final eft corresponds to the BWT acter moved to back) result	xam the s lly co . Ori ing ir	tep mp igir n di	es needed oute Shar all seque afferent s	l to nnce ence etrin	obton's e is r	tain the entrope otated as ma	e fir by. 1 l seq ny a xico	nal s) Th uen s ch grap	equence ne table tially (fir aracters phic ord
able Thes numl the le chara the s	0.1: CR explained by enter the tables summarizes there to the tables summarizes there is the table of tab	xam the s lly co T. Ori ing in quence to th	tepompigir dices	e. s needed oute Sha hal seque ifferent s are then Index of	d to nnce ence etrin n so	o obton's e is range, orteo is o	ain the entrope otated as made in learning the otated as made in learning the otated as the otated a	e fir by. 1 l seq ny a xico g (e	nal s) Th uen s ch grap	equence ne table tially (fir aracters phic ord the thi
able Thes numl the le chara the s The seque	0.1: CR explained by enter three tables summarizes there from which we will final eff corresponds to the BWT except moved to back) result sequence. The resulting sequence in corresponds the ence here in original order	xam the solly cooling in quenction the thick is a second	tepompigin dices	es needed bute Shanal seque afferent share then are then akes the	l to nnce ence trin n so th	o obton's e is rangs, orteo tis ofth p	ain the entrope otated as madering the control of t	e fir by. 1 l seq ny a xico g (e n in	nal s) Th luen s ch grap e.g.:	equence ne table tially (fi- aracters phic ord the thi- icograph
able Thes numl the le chara the s The seque order	0.1: CR explained by enter three tables summarizes there from which we will final eff corresponds to the BWT excter moved to back) result sequence. The resulting sequence here in original order to the control of the	xam the s ly co ly co ring in quenc to th "#"	ple mpigira di ces ne l' crre	es needed bute Sharal seque afferent sare then Index of akes the esponds	l to nnce trin n so th fif	o obton's e is ranges, orteo tis out the part of the part of the part of the outer	ain the entropotated as madin le rderinositio result	e fir by. 1 l seq ny a xico g (e n in of t	nal s) Th luen s ch grap e.g.: lex he H	equence ne table tially (fi- aracters phic ord the the icograph BWT, the
These numbers the sequence order is the	0.1: CR explained by ease three tables summarizes ber from which we will final eft corresponds to the BWT eacter moved to back) result sequence. The resulting sequence here in original order r). 2) The table in the center last character of previous	the solly cooling in queno the "#' ter co	ple mp igir n di ces ne l' ta	e. s needed bute Sharal seque afferent s are then Index of akes the esponds ces orde	l to nnce ttrir n so th fif to	o obton's e is rangs, orted the pthe	eain the entropotated as made in lear dering consition result	te firm by. 1 sequence of the	nal s) Th uen s ch grap e.g.: lex he H	equence ne table tially (firaracters phic ord the thin icograph 3WT, the table
Thes number the leaders the sequence order is the results the resu	o.1: CR explained by ease three tables summarizes ber from which we will final eft corresponds to the BWT eacter moved to back) result sequence. The resulting sequence here in original order r). 2) The table in the center last character of previous eight corresponds to the ap	xam the s lly co	teple igir n di ces ne l orre	es needed bute Shar hal seque afferent s are then Index of akes the esponds ces orde n of the	l to mnce trim n so th fif to MI	o obton's e is rugs, orteo is of the as of F	eain the entropotated as made in lear dering cosition result explain algoritic control of the co	e firm by. 1 sequence of the s	nal s) Th tuents ch grap e.g.: lex he H Sta	equence ne table tially (fi- aracters bhic ord the the ticograph BWT, the table rting from
able Thes numl the le chara the s The seque order is the	0.1: CR explained by ease three tables summarizes ber from which we will final eft corresponds to the BWT eacter moved to back) result sequence. The resulting sequence here in original order r). 2) The table in the center last character of previous	xam the s lly co	teple igir n di ces ne l orre	es needed bute Shar hal seque afferent s are then Index of akes the esponds ces orde n of the	l to mnce trim n so th fif to MI	o obton's e is rugs, orteo is of the as of F	eain the entropotated as made in lear dering cosition result explain algoritic control of the co	e firm by. 1 sequence of the s	nal s) Th tuents ch grap e.g.: lex he H Sta	equence ne table tially (fi aracters phic ord the the icograph BWT, the ne table rting from
able Thes numl the le chara the s The seque order is the the r a sec	o.1: CR explained by ease three tables summarizes ber from which we will final eft corresponds to the BWT eacter moved to back) result sequence. The resulting sequence here in original order r). 2) The table in the center last character of previous eight corresponds to the ap	xam the s lly co l. Ori ing in quence to th '"#' ter co sequence plica ned l	teple tepper tep	es needed oute Shar hal seque different s are then Index of akes the esponds ces orde n of the e "Char	l to mnce ttrir n so th fif to red M7	obton's e is rangs, sorted is continued the as a continued as a continued the action to the action the action that a continued	tain the entrope otated as mad in lear dering continuous tion result explain learning from the form of the explain learning from the explain learnin	te firm by. 1 sequence of the	nal s) Th luents ch grap e.g.: lex he H Sta nuc	equence ne table tially (fi aracters phic ord the the icograph BWT, the ne table arting from
able Thes numl the le chara the s The seque order is the the r a sec this e	o.1: CR explained by ease three tables summarizes ber from which we will final eft corresponds to the BWT acter moved to back) result sequence. The resulting search eric here in original order r. 2) The table in the centre last character of previous eight corresponds to the apquence of all character narrosse), the MTF will get the	xam the s lly co lly co co to the ter co sequence plica ned l	teple	es needed oute Shar all seque afferent share then Index of akes the esponds ces orden of the e "Char of the cu	l to mnce trim so th fif to MI li	obton's e is rangs, portection of the portection	tain the entrope otated as madering the control of	e fir by. 1 seq ny a xico g (ϵ n in of t hm. 3	nal s) Th uen s ch grap lex.g.: lex he H Sta nuc	equence ne table tially (fi aracters phic ord the th icograp BWT, the ne table rting fro leotides the BV
Thes number the seque order is the rate as each this control of the seque order is the rate as sec this control of the seque order is the rate of the seque of th	o.1: CR explained by enter three tables summarizes there from which we will final eff corresponds to the BWT eacter moved to back) result sequence. The resulting sequence here in original order r. 2) The table in the centre elast character of previous eight corresponds to the appunce of all character nances, the MTF will get the er case bold letter) in the	xam the s ly co ly co ring in quence to the sequence sequence plica med le inde	ple igir igir igir igir igir igir igir igi	es needed oute Shanal seque different so are then Index of akes the esponds ces orde on of the e "Char of the cu list".	l to nnce trir n so th fif to MI . li rred	obtoon's e is reasonable to obtoon's e is reasonable to obtoon the properties of the control of	tain the entropotated as made in lear result explaint (our function of the explaint control of the explaint (our function our	e fir yy. 1 l seq my a xico g (e n in of t n. 3 hm. our ide f	nal s) Th quent s ch ggrap e.g.: lex he F Sta nuc rom o, fo	equence table tially (fi aracters phic ord the thricograph BWT, the table rting from the BW r the new table rting from the BW r the new table requirements the BW r the new table requirements the BW r the new table requirements the same table requirements t
able Thes numl the le chara the s The seque order is the the r a sec this e (upp itera	o.1: CR explained by ease three tables summarizes ber from which we will final eft corresponds to the BWT acter moved to back) result sequence. The resulting search eric here in original order r. 2) The table in the centre last character of previous eight corresponds to the apquence of all character narrosse), the MTF will get the	xam the s lly co c	tep mpigirn di ne di ces ne di cer tion her ex conar.	es needed oute Shared seque afferent seque are then are then are the esponds ces orde nof the e"Char of the culist".	l to nnce trir n so th fif to red Mil. li In	obtoon's e is rengs, sorted is of the properties	tain the entropotated as made in lear result explaint (our function of the explaint control of the explaint (our function our	e fir yy. 1 l seq my a xico g (e n in of t n. 3 hm. our ide f	nal s) Th quent s ch ggrap e.g.: lex he F Sta nuc rom o, fo	equence table tially (fi aracters phic ord the thricograph BWT, the table rting from the BW r the new table rting from the BW r the new table requirements the BW r the new table requirements the BW r the new table requirements the same table requirements t

Finally, we compute the Shannon's entropy of the values obtained by the MTF that results in our CR (the CV is obtained by multiplying CR by the length

of the sequence).

length(s)

$$p_{(i)} = \frac{N_i}{length(s)} \tag{0.1}$$

 $E(s) = -\sum p_{(i)} \times log_4(p_{(i)})$

$$CR(s) = E(MTF(BWT(s))) \tag{0.3}$$

$$CV(s) = E(MTF(BWT(s))) \times length(s)$$
 (0.4)

0.1.2 Complexity in strings

Genomic sequences

Annotation of repetitive elements

Human texts

Complexity in windows

0.1.3 Simulations

0.2 Measuring dynamics of genetic species

0.2.1 Genomes

0.2.2 Mining of Genetic Elements

Repetitive Elements

Functional Elements

0.2.3 Randomization of genetic elements

Chr. name	Original length	Corrected length	Percentage le
1	249,240,621	225,200,000	90.35%
2	243,188,741	237,670,000	97.73%
3	197,961,181	194,230,000	98.12%
4	191,044,271	187,270,000	98.02%
5	180,901,928	177,090,000	97.89%
6	171,048,878	167,050,000	97.66%
7	159,128,663	154,640,000	97.18%
8	146,302,151	142,290,000	97.26%
9	141,151,937	120,130,000	85.11%
10	135,524,747	131,040,000	96.69%
11	134,946,516	130,310,000	96.56%
12	133,841,891	129,970,000	97.11%
13	115,109,733	95,500,000	82.96%
14	107,289,415	87,910,000	81.94%
15	102,521,389	81,520,000	79.52%
16	90,290,985	78,640,000	87.10%
17	81,195,208	77,700,000	95.70%
18	78,017,245	74,580,000	95.59%
19	59,118,983	55,460,000	93.81%
20	62,962,324	59,430,000	94.39%
21	48,119,895	35,100,000	72.94%
22	51,244,541	34,790,000	67.89%
X	155,260,558	150,230,000	96.76%
Y	59,033,288	22,520,000	38.15%
Table 0.2: Transformation of Chromosome size. Example of changes in estimation of chromosome length after removing regions with no GEs for Human chromosome 1.			

0.2.4 Ecolopy

0.2.5 Neutral Ecological models

Ewens sampling formula

$$\theta = 2J_M \nu \tag{0.5}$$

$$Pr\{S, n1, n2, \dots, n_S | \theta\} = \frac{J_M! \theta^S}{1^{\phi_1} 2^{\phi_2} \cdots J_M^{\phi_{J_M}} \phi_1! \phi_2! \cdots \phi_{J_M}! \prod_{k=1}^{J_M} (\theta + k - 1)}$$
(0.6)

$$\mathcal{L} = \frac{\theta^S}{\prod_{k=1}^{J_M} (\theta + k - 1)} \tag{0}$$

Etienne sampling formula

$$m = \frac{I}{I + J - 1} \tag{0.8}$$

$$P[D|\theta, m, J] = \frac{J!}{\prod_{i=1}^{S} n_i \prod_{j=1}^{J} \phi_J!} \frac{\theta^S}{(I)_J} \sum_{A=S}^{J} K(D, A) \frac{I^A}{(\theta)_A}$$
(0.9)

 $\prod_{i=1}^{S} \frac{\bar{s}(n_i, a_i)\bar{s}(a_i, 1)}{\bar{s}(n_i, 1)}$

 $\bar{s}(n_i,1)$

K(D,A) :=

 $\{a_1,...,a_s | \sum_{i=1}^{S} a_i = A\}^{i=1}$

(0.10)

$$S_{(n,m)} = S_{(n-1,m-1)} - (n-1) \times S_{(n-1,m)}$$
(0.11)

$$P[D|\theta, m, J] = \frac{J!}{\prod_{i=1}^{S} n_i \prod_{J=1}^{J} \Phi_J!} \frac{\theta^S}{(\theta)_J} \times \sum_{A=S}^{J} \left(K(D, A) \frac{(\theta)_J}{(\theta)_A} \frac{I^A}{(I)_J} \right)$$
(0.12)

0.2.6 Model optimization

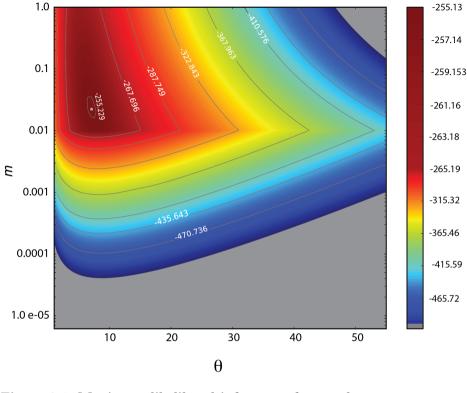
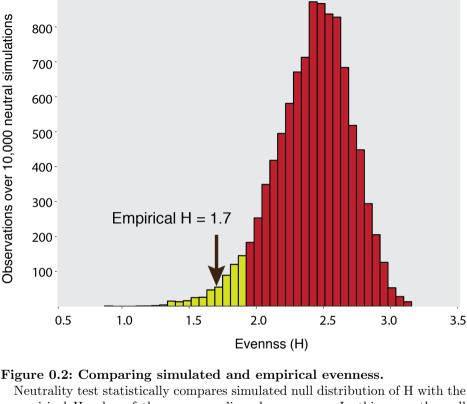


Figure 0.1: Maximum likelihood inference of neutral parameters. Log likelihood surface as a function of migration rate (m), and the fundamental

biodiversity number (θ) for D. rerio chromosome 19. Dark red color shows regions of the surface where parameters maximize the probability to explain abundances and diversity of genetic elements in the chromosome. Likelihood ratio tests favored Etienne in contrast to Ewens sampling formula to explain the observed data in the chromosome.

0.2.7 Model testing

0.2.8 Test for neutrality



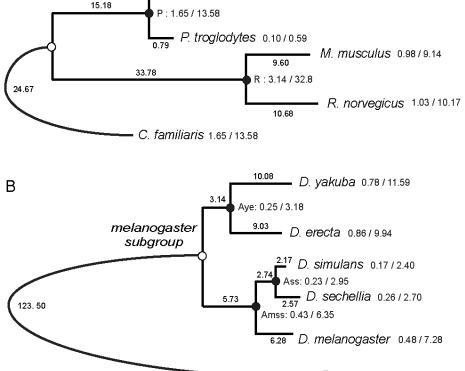
900

empirical H value of the corresponding chromosome. In this case, the null distribution of H values was derived from 10,000 neutral simulations of A. gambiae chromosome 2L, with neutral parameters (θ and m) optimized by ML using Etienne sampling formula. Yellow and red bars display 5% and 95% of the simulated neutral data, respectively. Although in this case neutrality was rejected (p= 0.01), posterior correction by multiple testing favored the null

neutral hypothesis (q = 0.21).

0.3 Detection of selective pressure at molecular level

0.3.1 Orthology prediction



H. sapiens 0.08/0.56

D. ananassae 7.90 / 151.07 Figure 0.3: Mammals and *melanogaster* group phylogeny.

Numbers on internal and external nodes represent the median number of nonsynonymous and synonymous substitutions per codon (dN/dS) estimated from all the coding sequences compared in mammal (A) and Drosophila (B) genomes. Branch lengths and rates were multiplied by 100. Ancestral estimation

of parameters was done in primates (P), rodents (R), D. yakuba and D. erecta (Aye), D. simulans and D. sechellia (Ass), and D. melanogaster, D. simulans

and D. sechellia (Amss). C. familiaris and D. ananassae were chosen as out-

group species in the corresponding tree.

0.3.2 Alignments refinement and filters

0.3.3 Evolutionary analysis

0.3.4 GSSA, evolutionary and statistical simulations

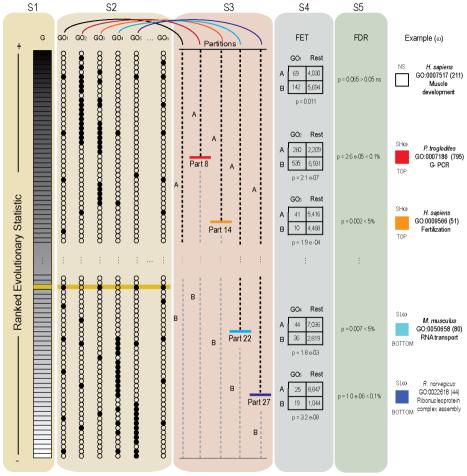


Figure 0.5:

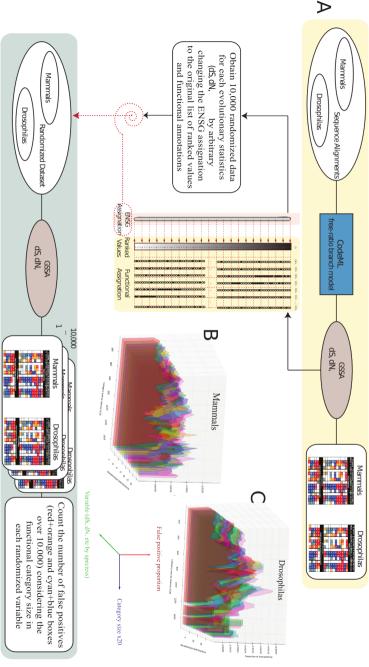


Figure 0.7:

0 11 L) and and مال ا on experiment...

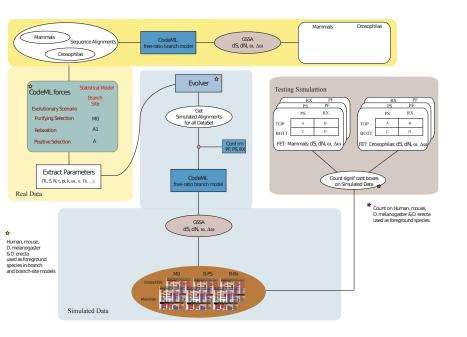


Figure 0.8: Evolutionary and statistical simulation of GSSA.

The pipeline shows the steps taken along three different spaces of analysis, the real data, the simulated data and the testing block. See Supplementary Results for a complete explanation of methods and results.

	PS		RX		PF	
	# PSG	# RXG	# PSG	# RXG	# PSG	# RXG
Homo sapiens	658	1640	11	1939	0	1
Mus musculus	1500	954	14	1565	1	0
D. melanogaster	736	630	25	1104	0	0
D. erecta	778	1292	26	1713	2	1
Table 0.3: Number evolutionary scen		and relax	ed genes ((RXG) in (each of the	e simulated

PS	_	92.50%	98.50%
RX	91.10%		99.00%
PF	88.90%	90.60%	—

RX

PS

Table 0.4: Proportion of significant functional categories that are still significant (identical signs of odd-ratios) under a different evolutionary scenario.