

SENCA: A codon substitution model to better estimate evolutionary processes

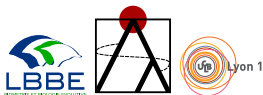
—

Fanny Pouyet

Marc Bailly-Bechet, Dominique Mouchiroud
and Laurent Guéguen

—

LBBE - UCB Lyon - UMR 5558
June 2015, Porquerolles, France.

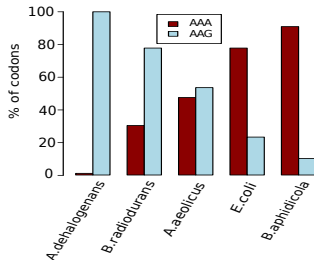


Codon Usage Bias - CUB

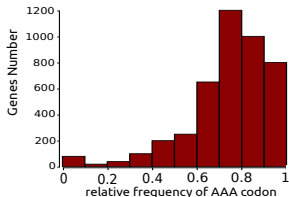
	U	C	A	G
U	UUU } Phe UUC } UUA } Leu UUG }	UCU } Ser UCC } UCA } UCG }	UAU } Tyr UAC } UAA } Stop UAG }	UGU } Cys UGC } UGA } Stop UGG } Trp
C	CUU } Leu CUC } CUA } CUG }	CCU } Pro CCC } CCA } CCG }	CAU } His CAC } CAA } Gln CAG }	CGU } Arg CGC } CGA } CGG }
A	AUU } Ile AUC } AUA } AUG } Met	ACU } Thr ACC } ACA } ACG }	AAU } Asn AAC } AAA } Lys AAG }	AGU } Ser AGC } AGA } Arg AGG }
G	GUU } Val GUC } GUA } GUG }	GCU } Ala GCC } GCA } GCG }	GAU } Asp GAC } GAA } Glu GAG }	GGU } Gly GGC } GGA } GGG }

61 sense codons, 20 amino acids

Lysine codon usage in bacteria



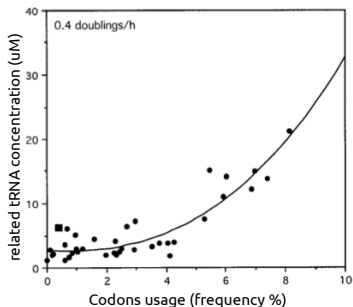
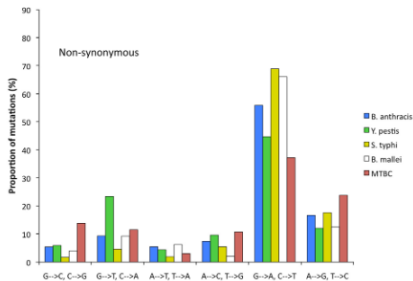
E.coli Lysine Usage



CUB

How to measure it?

Origins of the CUB: mutational vs. selective explanations.



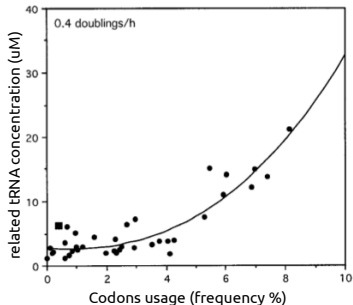
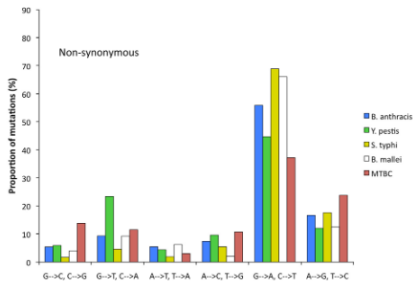
Objectives: Disentangling the two hypotheses

Modelling evolutive processes explaining the observed CUB: at the nucleotidic (N), the codons (C) and the AA layers (A).

CUB

How to measure it?

Origins of the CUB: mutational vs. selective explanations.

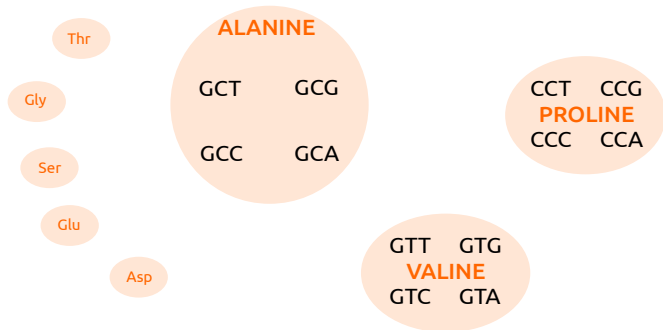


Objectives: Disentangling the two hypotheses

Modelling evolutive processes explaining the observed CUB: at the nucleotidic (N), the codons (C) and the AA layers (A).

SENCA: Sites Evolution of Nucleotides, Codons and Amino-acids

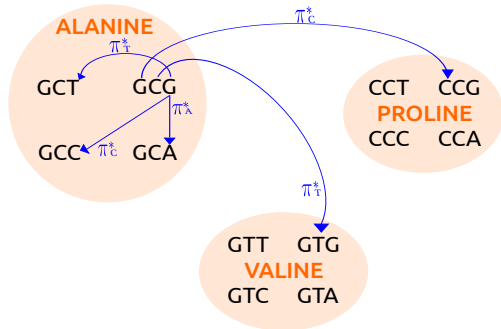
Let's have an example



Inspired by Yang and Nielsen 2008
(FMutSel).

SENCA: Sites Evolution of Nucleotides, Codons and Amino-acids

N layer

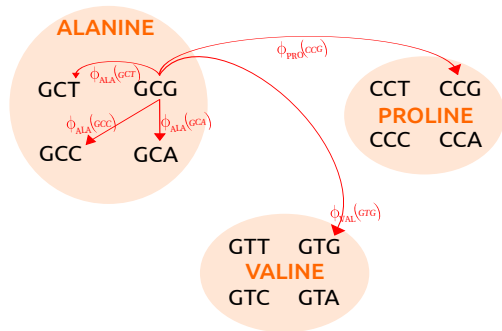


κ : transition/transversion

π_n^* : equilibrium frequency of nucleotide $n \in [A, C, G, T]$.

SENCA: Sites Evolution of Nucleotides, Codons and Amino-acids

C layer

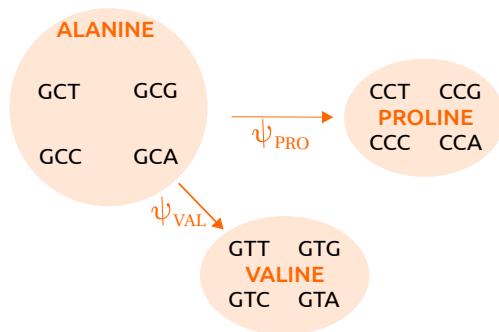


$\phi_{AA}(i)$: codon i preference – intra-AA.

$$\sum_{I \in \text{codons}(AA)} \phi_{AA}(I) = 1$$

SENCA: Sites Evolution of Nucleotides, Codons and Amino-acids

A layer

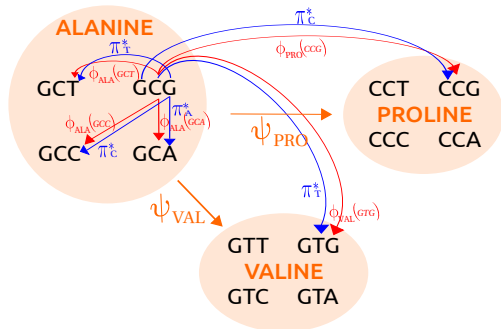


ω : non-synonymous/synonymous substitution rates

ψ_{AA} : preference AA. $\sum_{AA \in \text{amino acids}} \psi(AA) = 1$

SENCA: Sites Evolution of Nucleotides, Codons and Amino-acids

Inspiration and structure



Implemented in Bio++

<http://biopp.univ-montp2.fr/>.

SENCA

Formulas

Instantaneous evolutionary rate from codon i to j :

$$q_{ij} = \begin{cases} 0 & \text{if 2 or 3 different positions,} \\ \pi_{jk}^* \kappa f(x_i, x_j) & \text{synonymous transition,} \\ \pi_{jk}^* \kappa \omega f(x_i, x_j) & \text{non-synonymous transition,} \\ \pi_{jk}^* f(x_i, x_j) & \text{synonymous transversion,} \\ \pi_{jk}^* \omega f(x_i, x_j) & \text{non-synonymous transversion.} \end{cases}$$

with:

$$f(x_i, x_j) = \frac{-\log\left(\frac{x_i}{x_j}\right)}{1 - \frac{x_i}{x_j}}$$

and:

$$x_i = n_{aa} \phi_{aa}(i) \psi_{aa_i}$$

Data

21 pathogeneus bacteria and archea

From Lassalle et al., *PLoS Genet.* 2015 \Rightarrow concatenates of approx. 100 genes by increasing ENC (core genome, non-recombinant).

Dataset	Taxon Name	Nb. of strains	Nb of concatenates	Mean GC %
brucella	<i>Brucella</i> spp.	9	4	58.8
francis	<i>Francisella tularensis</i>	8	3	33.8
mycobacterium	<i>Mycobacterium tuberculosis</i> complex	7	1	66.1
burk_mal	<i>Burkholderia Pseudomallei</i> group	9	6	68.7
yersinia	<i>Yersinia pestis</i>	11	6	49.3
chlamydia	<i>Chlamydia trachomatis</i>	13	4	41.8
sulfo	<i>Sulfolobus</i> spp.	8	4	35.4
burk_ceno	<i>Burkholderia cenocepacia</i> complex	8	8	68.2
staph	<i>Staphylococcus aureus</i>	15	5	34.2
bifido	<i>Bifidobacterium longum</i>	6	3	61.9
acineto	<i>Acinetobacter</i> spp.	6	5	40.8
campylo	<i>Campylobacter jejunii</i>	6	4	31.6
clostridium	<i>Clostridium botulinum</i>	8	5	29.6
strep_py	<i>Streptococcus pyogenes</i>	12	3	39.6
strep_pneu	<i>Streptococcus pneumoniae</i>	13	3	42.0
salmonella	<i>Salmonella enterica</i>	14	6	54.6
listeria	<i>Listeria</i> spp.	8	3	38.8
B_anthraxis	<i>Bacillus anthracis/aureus</i> group	17	3	37.0
nesseiria	<i>Nesseiria meningitidis</i>	8	2	55.3
escherichia	<i>Escherichia coli</i>	35	1	53.3
helicobacter	<i>Helicobacter pylori</i>	14	1	40.4

Implementation

- Non-stationary and homogeneous run,
- Intra-species runs: amino-acids preferences stationary,
- Optimisation by max. likelihood: 109 parameters.

Validation of the model:

Comparisons with AIC criterium to standard codons model YN98xF61.

SENCA is better in 68 over 78 concats. Exceptions are:

- 2 concats. of *Brucella* spp. and *Burkholderia peusomallei*,
- 1 of *Salmonella enterica* and *Yersinia pestis*,
- all of *Sulfolobus* spp.

Implementation

- Non-stationary and homogeneous run,
- Intra-species runs: amino-acids preferences stationary,
- Optimisation by max. likelihood: 109 parameters.

Validation of the model:

Comparisons with AIC criterium to standard codons model YN98xF61.

SENCA is better in 68 over 78 concats. Exceptions are:

- 2 concats. of *Brucella* spp. and *Burkholderia peusomallei*,
- 1 of *Salmonella enterica* and *Yersinia pestis*,
- all of *Sulfolobus* spp.

What about the CUB?

- To disentangle between the 2 hypotheses: compute the CUB by comparison with expected if uniform \rightarrow ENC index.
- Effective Number of Codons – ENC – varies between 61 and 20.

	U	C	A	G
U	UUU } Phe UUC } UUA } UUG } Leu	UCU } Ser UCC } UCA } UCG } Stop	UAU } Tyr UAC } UAA } UAG } Stop	UGU } Cys UGC } UGA } Stop UGG } Trp
C	CUU } CUC } Leu CUA } CUG } Leu	CCU } Pro CCC } CCA } CCG } Leu	CAU } His CAC } CAA } Gln CAG } Pro	CGU } CGC } Arg CGA } CGG } Arg
A	AUU } Ile AUC } AUA } AUG } Met	ACU } Thr ACC } ACA } ACG } Thr	AAU } Asn AAC } AAA } Lys AAG } Asn	AGU } Ser AGC } AGA } Arg AGG } Arg
G	GUU } Val GUC } GUA } GUG } Val	GCU } Ala GCC } GCA } GCG } Ala	GAU } Asp GAC } GAA } Glu GAG } Asp	GGU } Gly GGC } GGA } GGG } Gly

$ENC = 61$

	U	C	A	G
U	UUU } Phe UUC } UUA } UUG } Leu	UCU } Ser UCC } UCA } UCG } Stop	UAU } Tyr UAC } UAA } UAG } Stop	UGU } Cys UGC } UGA } Stop UGG } Trp
C	CUU } CUC } Leu CUA } CUG } Leu	CCU } Pro CCC } CCA } CCG } Leu	CAU } His CAC } CAA } Gln CAG } Pro	CGU } CGC } Arg CGA } CGG } Arg
A	AUU } Ile AUC } AUA } AUG } Met	ACU } Thr ACC } ACA } ACG } Thr	AAU } Asn AAC } AAA } Lys AAG } Asn	AGU } Ser AGC } AGA } Arg AGG } Arg
G	GUU } Val GUC } GUA } GUG } Val	GCU } Ala GCC } GCA } GCG } Ala	GAU } Asp GAC } GAA } Glu GAG } Asp	GGU } Gly GGC } GGA } GGG } Gly

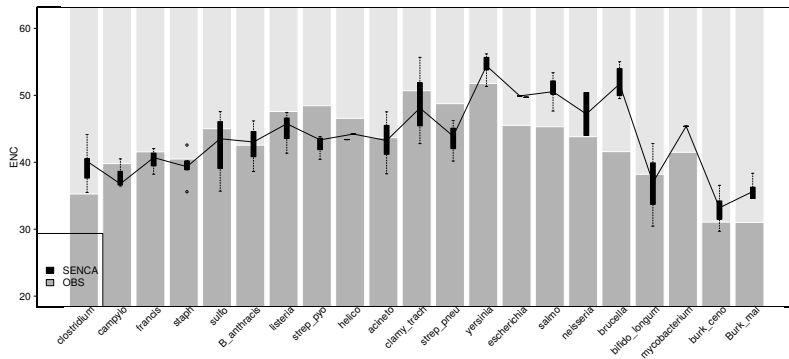
$ENC = 32$

	U	C	A	G
U	UUU } Phe UUC } UUA } UUG } Leu	UCU } Ser UCC } UCA } UCG } Stop	UAU } Tyr UAC } UAA } UAG } Stop	UGU } Cys UGC } UGA } Stop UGG } Trp
C	CUU } CUC } Leu CUA } CUG } Leu	CCU } Pro CCC } CCA } CCG } Leu	CAU } His CAC } CAA } Gln CAG } Pro	CGU } CGC } Arg CGA } CGG } Arg
A	AUU } Ile AUC } AUA } AUG } Met	ACU } Thr ACC } ACA } ACG } Thr	AAU } Asn AAC } AAA } Lys AAG } Asn	AGU } Ser AGC } AGA } Arg AGG } Arg
G	GUU } Val GUC } GUA } GUG } Val	GCU } Ala GCC } GCA } GCG } Ala	GAU } Asp GAC } GAA } Glu GAG } Asp	GGU } Gly GGC } GGA } GGG } Gly

$ENC = 20$

Strength of CUB

ENC index



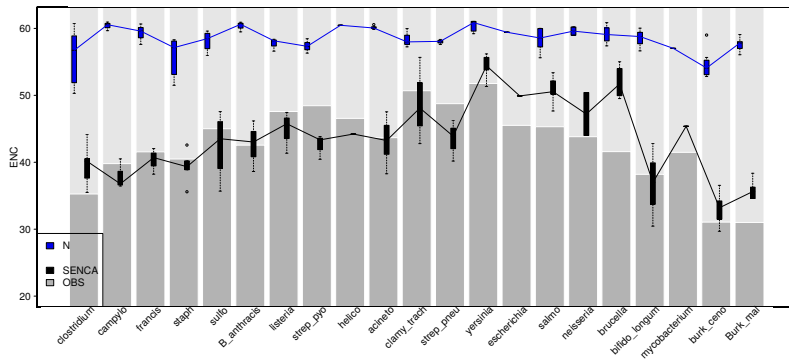
At equilibrium, frequency of codon i is proportional to:

$$f^*(i) \propto \prod_{k=1}^3 \pi_{ik}^* \times \phi_{aa_i}(i) n_{aa_i} \times \psi_{aa_i}$$

$$\rightarrow ENC_{SENCA}^* \approx ENC_{OBS}$$

Species ordered by incr. GC_{obs}

ENC index



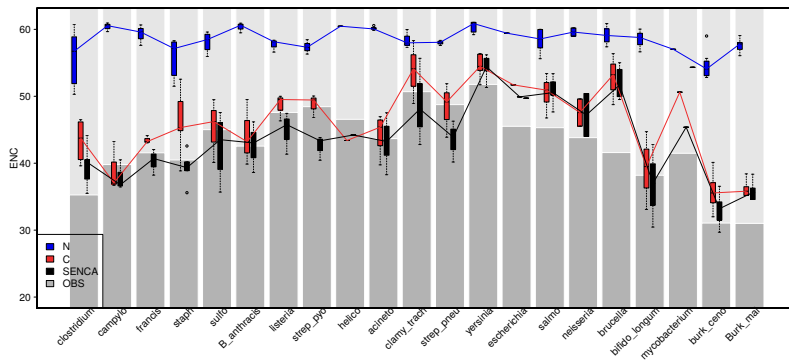
Compute ENC_N^* :

$$f^*(i) \propto \prod_{k=1}^3 \pi_{i_k}^* \times 1 \times 1$$

→ $ENC_N^* \approx$ uniform CUB,

Species ordered by incr. GC_{Obs}

ENC index



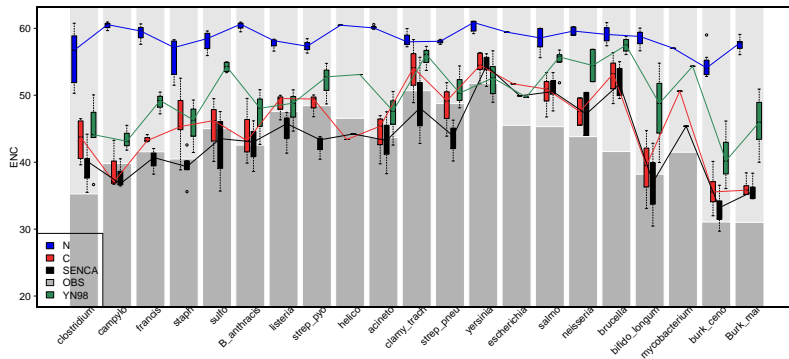
Compute ENC_C^* :

$$f^*(i) \propto 1 \times \phi_{aa_i}(i) n_{aa_i} \times 1$$

$$\rightarrow ENC_C^* \approx ENC_{SENCA}^*$$

Species ordered by incr. GC_{obs}

ENC index

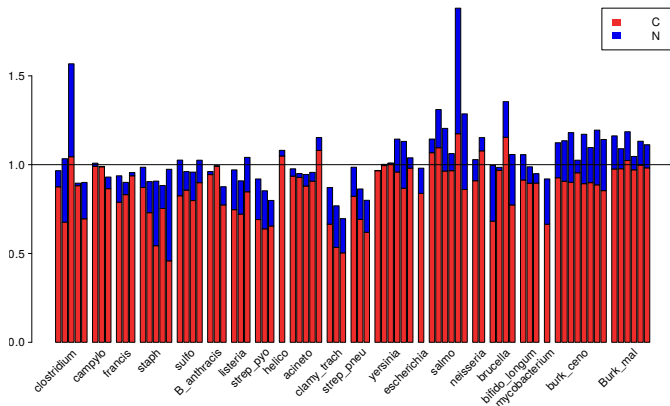


- $ENC_{SENCA}^* \approx ENC_{OBS} < ENC_{YN98}^*$,
- Is ENC_{SENCA}^* mostly driven by ENC_C^* ? Quantifying.

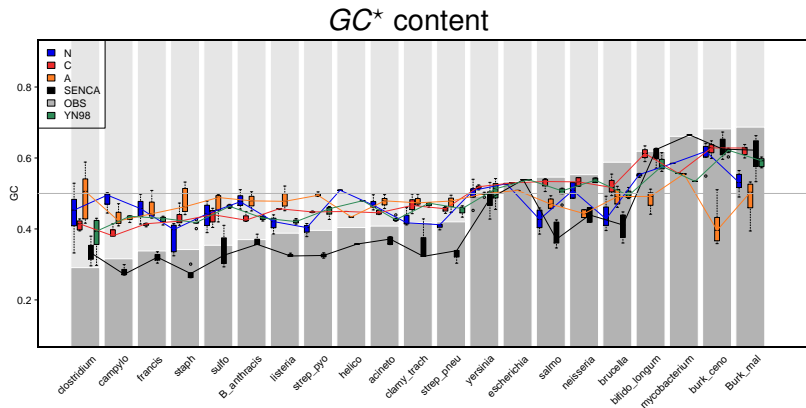
Quantification

$$dENC_{SENCA}^* \approx dENC_C^* + dENC_N^*$$

- CUB measured as distance to uniform usage: $dENC^* = 61 - ENC^*$,
- High correlation ($R = 0.95$, $p\text{-value} < 10e-16$, Pearson correlation test),
- Relative importance of C and N : $\frac{dENC_C^*}{dENC_{SENCA}^*} > \frac{dENC_N^*}{dENC_{SENCA}^*}$.

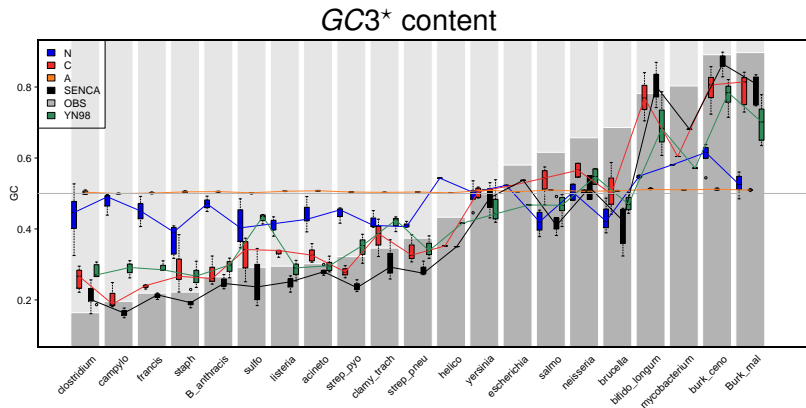


What about genome composition?



- $GC^*_{SENCA} < GC_{obs}$
- SENCA combined effects of the 3 layers → AT enrichment at equilibrium.

What about genome composition?

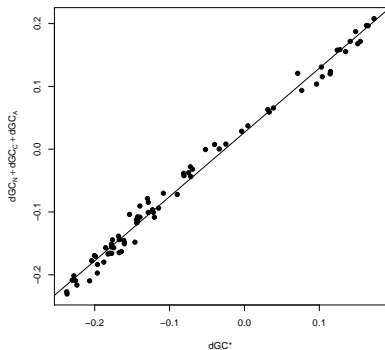


- $GC3^*_{SENCA}$ more biased than $GC3^*_{YN98}$: in agreement with $GC3_{obs}$,

⇒ Which layer has the most important impact on GC bias?

Quantification

$$dGC^* \approx dGC_N^* + dGC_C^* + dGC_A^*$$

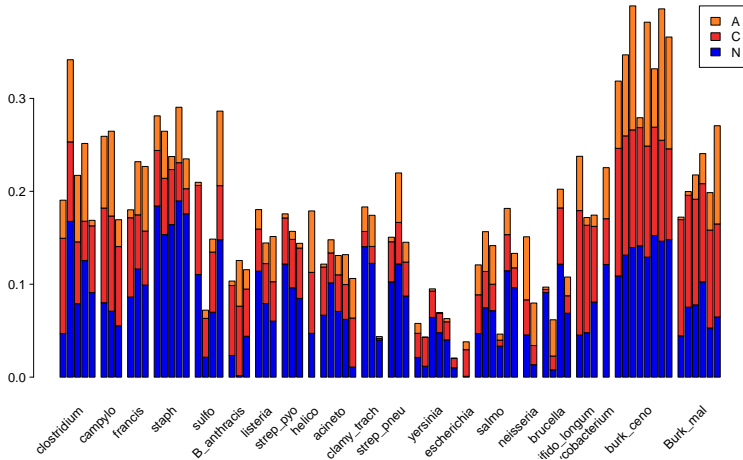


GC^* bias measured as $dGC^* = 0.5 - GC^*$ (distance to uniform composition).
High correlation ($R=0.99$, $p\text{-value} < 10e-16$, Pearson correlation test).

Slope=1.04, intersect fixed to 0.

Quantification

$$dGC^* \approx dGC_N^* + dGC_C^* + dGC_A^*$$



N, **C** and **A** impact on GC^* .

Conclusion

- **Methodologically:**
 - Non-stationary model
 - Multilayered model
 - New statistical tools: GC and ENC of layers

- **Biologically:**
 - Bias towards AT at equilibrium
 - Importance of N , C and A for genomic bias.
 - CUB mostly due to C .

Perspectives

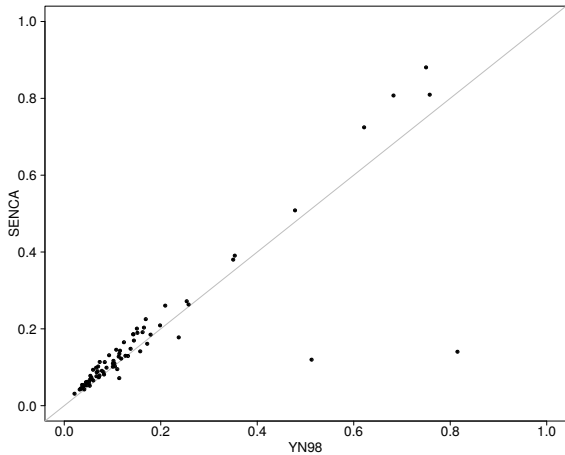
- **Extension of the model:**
 - Gene expression factor : measure of CUB intensity within a genome
- **Biological questions:**
 - Influence on ω estimation,
 - HIV – Human interaction and consequences on HIV CUB,
 - Can we detect recombination patterns: BGC?

Thank you for your attention!



- **Acknowledgments:**
Laurent Guéguen, Marc Bailly-Bechet, Dominique Mouchiroud, Florent Lassalle, Vincent Daubin

Omega comparisons



SENCA: Site Evolution of Nucleotides, Codons and Amino-acids

Implementation and parameters

SENCA with 65 parameters:

π_{jk}^* equil. frequency of mutational process of j_k

κ transition/transversion

$\phi_{aa}(i)$ codon preference. For each AA, we have: $\sum_i \phi_{aa}(i) = 1$

n_{aa_i} degenerescence of the AA

$\omega = \frac{dN}{dS}$ mutation rate of non-synonymous dN and synonymous dS

ψ_{aa} AA preference. We have: $\sum_{aa} \psi_{aa} = 1$

In Bio++ (<http://biopp.univ-montp2.fr/>). Optimisation by max. likelihood. Homogeneous, non-stationary analysis (ψ_{aa} stat. because intra-species analyses)