

# 1 SkewDB, a comprehensive database of GC and 10 2 other skews for over 30,000 chromosomes and 3 plasmids

4 Bert Hubert\*

5 \*corresponding author: Bert Hubert (bert@hubertnet.nl)

## 6 ABSTRACT

GC skew denotes the relative excess of G nucleotides over C nucleotides on the leading versus the lagging replication strand of eubacteria. While the effect is small, typically around 2.5%, it is robust and pervasive. GC skew and the analogous TA skew are a localized deviation from Chargaff's second parity rule, which states that G and C, and T and A occur with (mostly) equal frequency even within a strand.

7 Different bacterial phyla show different kinds of skew, and differing relations between TA and GC skew.

This article introduces an open access database (<https://skewdb.org>) of GC and 10 other skews for over 30,000 chromosomes and plasmids. Further details like codon bias, strand bias, strand lengths and taxonomic data are also included. The *SkewDB* can be used to generate or verify hypotheses. Since the origins of both the second parity rule and GC skew itself are not yet satisfactorily explained, such a database may enhance our understanding of prokaryotic DNA.

## 8 Background & Summary

9 The phenomenon of GC skew<sup>1-3</sup> is tantalizing because it enables a simple numerical analysis that accurately predicts the loci of  
10 both the origin and terminus of replication in most bacteria and some archaea<sup>4,5</sup>.

11 Bacterial DNA is typically replicated simultaneously on both strands, starting at the origin of replication<sup>6</sup>. Both replication  
12 forks travel in the 5' to 3' direction, but given that the replichores are on opposite strands, topologically they are traveling in  
13 opposite directions. This continues until the forks meet again at the terminus. This means that roughly one half of every strand  
14 is replicated in the opposite direction of the other half. The forward direction is called the leading strand. Many plasmids also  
15 replicate in this way<sup>7</sup>.

16 The excess of G over C on the leading strand is in itself only remarkable because of Chargaff's somewhat mysterious  
17 second parity rule<sup>8</sup>, which states that within a DNA strand, there are nearly equal numbers of G's and C's, and similarly, T's  
18 and A's. This rule does not directly follow from the first parity rule, which is a simple statement of base pairing rules.

19 Depending on who is asked, Chargaff's second parity rule is so trivial that it needs no explanation, or it requires complex  
20 mathematics and entropy considerations to explain its existence<sup>9</sup>.

21 The origins of GC skew are still being debated. The leading and lagging strands of circular bacterial chromosomes are  
22 replicated very differently; it is at least plausible that this leads to different mutational biases. In addition, there are selection  
23 biases that are theorized to be involved<sup>10</sup>. No single mechanism may be exclusively responsible.

24 This article does not attempt to explain or further mystify<sup>11</sup> the second parity rule or GC skew, but it may be that the  
25 contents of the *SkewDB* can contribute to our further understanding.

26 The *SkewDB* also contains some hard to explain data on many chromosomes. These include highly asymmetric skew, but  
27 also very disparate strand lengths. Conversely, the *SkewDB* confirms earlier work on skews in the Firmicute phylum<sup>12</sup>, and also  
28 expands on these earlier findings.

29 GC skew has often been investigated by looking at windows of DNA of a certain size. GC skew is computed as  
30  $(G - C)/(G + C)$  in a window of  $N$  bases, where  $G$  is the number of guanines and  $C$  the number of cytosines in that window. It  
31 has been found that the choice of window size impacts the results of the analysis. The *SkewDB* is therefore based exclusively  
32 on cumulative skew<sup>13</sup>, which sidesteps window size issues. For example, the sequence GGGCCC has a cumulative GC skew  
33 of zero, and as we progress through the sequence, this skew takes on values 1, 2, 3, 2, 1, 0. If the window size  $N$  is 6, the  
34 non-cumulative skew is also 0.

35 The result of a cumulative GC skew analysis is shown in figure 1. The analysis software fits a linear model on the skews,  
36 where it also compensates for chromosome files sequenced in the non-canonical direction, or where the origin of replication is  
37 not at the start of the file.

38 GC skew analysis is not new. As noted below, the DoriC database for example contains related data that is more precise for  
39 its stated purpose (finding the Origin of replication). The SkewIT database<sup>4</sup> similarly provides a metric of skew, and also comes  
40 with an online analysis tool.

41 Other work, like the Comparative Genometrics Database<sup>14</sup> and the Z Curve Database<sup>15</sup> has been foundational, but by dint  
42 of their age lack an analysis of the tens of thousands of DNA sequences that have become available since the initial availability  
43 of these databases.

44 *SkewDB* is funded to be updated monthly with the latest sequences from NCBI until 2026.

45 Other software that calculates GC skew is available, like for example GraphDNA<sup>16</sup>, GC Skewing<sup>17</sup> and GenSkew. The  
46 *SkewDB* delivers far more metrics however, also because it involves annotation data in its calculations. For ease of use, *SkewDB*  
47 is made available as a ready to use database, as well as in software form that reproduces this database exactly.

## 48 Methods

49 The *SkewDB* analysis relies exclusively on the tens of thousands of FASTA and GFF3 files available through the NCBI  
50 download service, which covers both GenBank and RefSeq. The database includes bacteria, archaea and their plasmids.

51 Furthermore, to ease analysis, the NCBI Taxonomy database is sourced and merged so output data can quickly be related to  
52 (super)phyla or specific species.

53 No other data is used, which greatly simplifies processing. Data is read directly in the compressed format provided by  
54 NCBI. All results are emitted as standard CSV files.

55 In the first step of the analysis, for each organism the FASTA sequence and the GFF3 annotation file are parsed. Every  
56 chromosome in the FASTA file is traversed from beginning to end, while a running total is kept for cumulative GC and TA  
57 skew. In addition, within protein coding genes, such totals are also kept separately for these skews on the first, second and third  
58 codon position. Furthermore, separate totals are kept for regions which do not code for proteins.

59 In addition, to enable strand bias measurements, a cumulative count is maintained of nucleotides that are part of a positive  
60 or negative sense gene. The counter is increased for positive sense nucleotides, decreased for negative sense nucleotides, and  
61 left alone for non-genic regions. A separate counter is kept for non-genic nucleotides.

62 Finally, G and C nucleotides are counted, regardless of if they are part of a gene or not.

63 These running totals are emitted at 4096 nucleotide intervals, a resolution suitable for determining skews and shifts.

64 In addition, one line summaries are stored for each chromosome. These lines include the RefSeq identifier of the  
65 chromosome, the full name mentioned in the FASTA file, plus counts of A, C, G and T nucleotides. Finally five levels of  
66 taxonomic data are stored.

67 Chromosomes and plasmids of fewer than 100 thousand nucleotides are ignored, as these are too noisy to model faithfully.  
68 Plasmids are clearly marked in the database, enabling researchers to focus on chromosomes if so desired.

## 69 Fitting

70 Once the genomes have been summarised at 4096-nucleotide resolution, the skews are fitted to a simple model.

71 The fits are based on four parameters, as shown in figure 1. `alpha1` and `alpha2` denote the relative excess of G over C  
72 on the leading and lagging strands. If `alpha1` is 0.046, this means that for every 1000 nucleotides on the leading strand, the  
73 cumulative count of G excess increases by 46.

74 The third parameter is `div` and it describes how the chromosome is divided over leading and lagging strands. If this number  
75 is 0.557, the leading replication strand is modeled to make up 55.7% of the chromosome.

76 The final parameter is `shift` (the dotted vertical line), and denotes the offset of the origin of replication compared to the  
77 DNA FASTA file. This parameter has no biological meaning of itself, and is an artifact of the DNA assembly process.

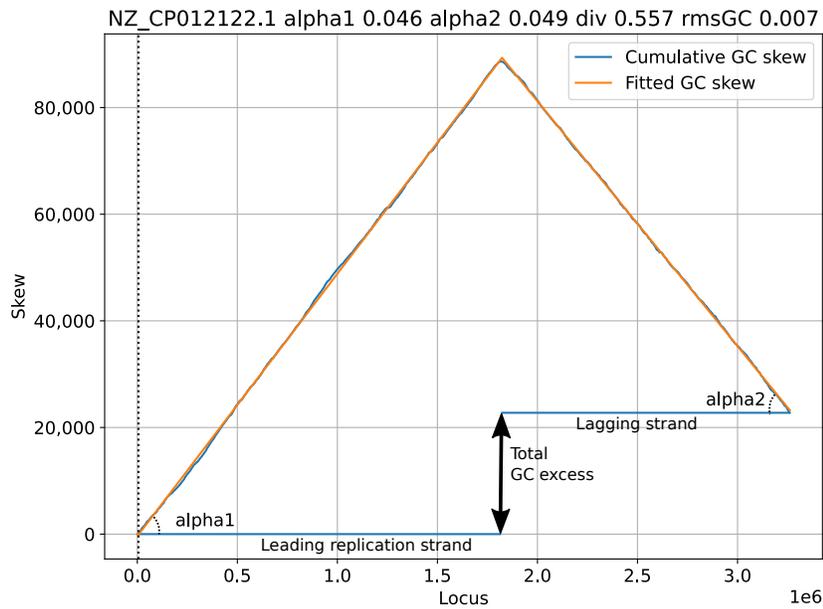
78 The goodness-of-fit number consists of the root mean squared error of the fit, divided by the absolute mean skew. This latter  
79 correction is made to not penalize good fits for bacteria showing significant skew.

80 GC skew tends to be defined very strongly, and it is therefore used to pick the `div` and `shift` parameters of the DNA  
81 sequence, which are then kept as a fixed constraint for all the other skews, which might not be present as clearly.

82 The fitting process itself is a downhill simplex method optimization<sup>18</sup> over the four dimensions, seeded with the average  
83 observed skew over the whole genome, and assuming there is no shift, and that the leading and lagging strands are evenly  
84 distributed. To ensure that the globally optimum fit is (very likely) achieved, ten optimization attempts are made from different  
85 starting points. This fitting process is remarkably robust in the sense that even significant changes in parameters or fitting  
86 strategies cause no appreciable change in the results.

87 For every chromosome and plasmid the model parameters are stored, plus the adjusted root mean squared error.

88 Both for quality assurance and ease of plotting, individual CSV files are generated for each chromosome, again at 4096  
89 nucleotide resolution, but this time containing both the actual counts of skews as well as the fitted result.



**Figure 1.** Sample graph showing *SkewDB* data for *Lactiplantibacillus plantarum* strain LZ95 chromosome

## 90 Some sample findings

91 To popularize the database, an online viewer is available on <https://skewdb.org/view>. While this article makes no independent  
 92 claims to new biological discoveries, the following sections show some results gathered from a brief study of the database.  
 93 Some of these observations may be of interest for other researchers.

## 94 GC and TA skews

95 Most bacteria show concordant GC and TA skew, with an excess of G correlating with an excess of T. This does not need  
 96 to be the case however. Figure 2 is a scatterplot that shows the correlation between the skews for various major superphyla.  
 97 Firmicutes (part of the Terrabacteria group) show clearly discordant skews.

## 98 Firmicute prediction

99 In many bacteria, genes tend to concentrate on the leading replication strand. If the codon bias of genes is such that they are  
 100 relatively rich in one nucleotide, the “strand bias” may itself give rise to GC or TA bias. Or in other words, if genes contain a  
 101 lot of G’s and they huddle on the leading strand, that strand will show GC skew. As an hypothesis, we can plot the observed GC  
 102 and TA skews for all Firmicutes for which we have data.

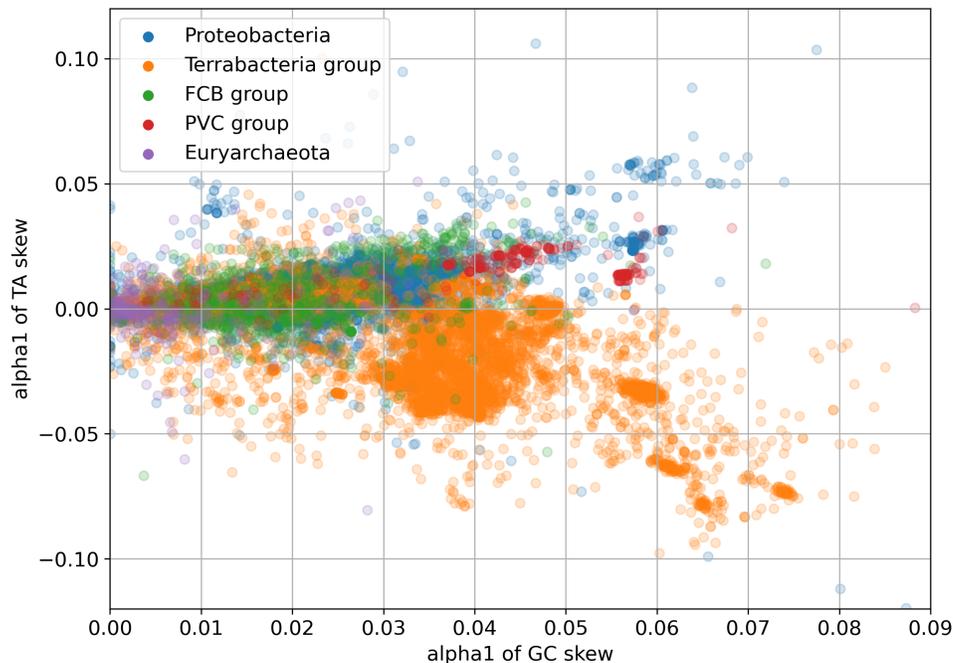
103 Mathematically the relation between the codon bias, strand bias and predicted GC skew turns out to be a simple multiplication.  
 104 In figure 3, the x-axis represents this multiplication. The y-axis represents the GC and TA skew ratio.

105 It can clearly be seen that both GC and TA skew correlate strongly with the codon/strand bias product. TA skew goes to  
 106 zero with the two biases, but GC skew appears to persist even in the absence of such biases.

107 Figure 4 shows the situation within an individual chromosome (*C. difficile*), based on overlapping 40960-nucleotide  
 108 segments. On the x-axis we find the strand bias for these segments, running from entirely negative sense genes to entirely  
 109 positive sense genes. The skew is meanwhile plotted on the y-axis, and here too we see that TA skew goes to zero in the absence  
 110 of strand bias, while GC skew persists and clearly has an independent strand-based component.

## 111 Asymmetric skew

112 The vast majority of chromosomes show similar skews on the leading and the lagging replication strands, something that  
 113 makes sense given the pairing rules. There are however many chromosomes that have very asymmetric skews, with one strand  
 114 sometimes showing no skew at all. In figure 5 four chromosomes are shown that exhibit such behavior. The *SkewDB* lists  
 115 around 250 chromosomes where one strand has a GC skew at least 3 times bigger/smaller than the other one.



**Figure 2.** Scatter graph of 25,000 chromosomes by superphylum, GC skew versus TA skew

### 116 **Asymmetric strands**

117 Bacteria tend to have very equally sized replication strands, which is also an optimum for the duration of replication. It is  
 118 therefore interesting to observe that GC skew analysis finds many chromosomes where one strand is four times larger than the  
 119 other strand. In figure 6 four such chromosomes are shown. The *SkewDB* lists around 100 chromosomes where one strand is at  
 120 least three times as large as the other strand.

### 121 **Anomalies**

122 In many ways, GC skew is like a forensic record of the historical developments in a chromosome. Horizontal gene transfer,  
 123 inversions, integration of plasmids, excisions can all leave traces. In addition, DNA sequencing or assembly artifacts will also  
 124 reliably show up in GC graphs, as elucidated with examples in<sup>4</sup>.

125 Figure 7 shows GC and TA skews for *Salmonella enterica subsp. enterica serovar Concord* strain AR-0407 (NZ\_CP044177.1),  
 126 and many things could be going on here. The peaks might correspond to multiple origins of replication, but might also indicate  
 127 inversions or DNA assembly problems.

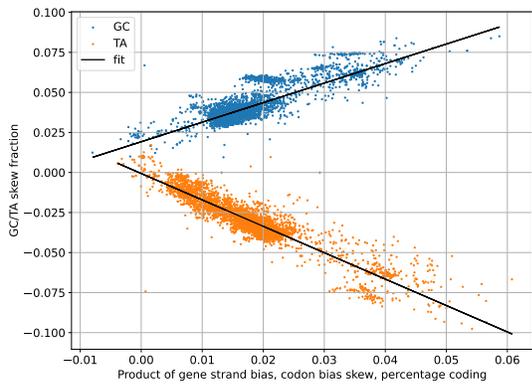
### 128 **Data Records**

129 The *SkewDB* is available through <https://skewdb.org>, where it is frequently (& automatically) refreshed. A snapshot of the data  
 130 has also been deposited on Dryad<sup>19</sup>.

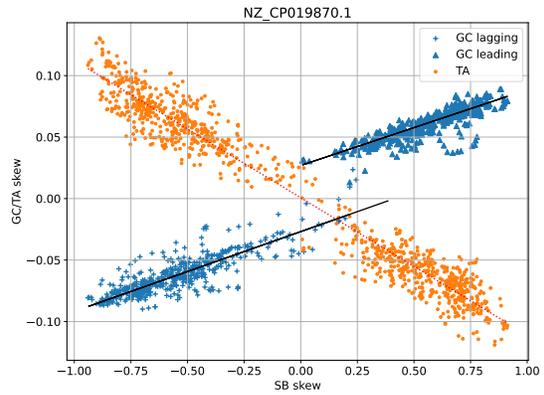
131 The *SkewDB* consists of several CSV files: *skplot.csv*, *results.csv*, *genomes.csv* and *codongc.csv*. In addition, for each  
 132 chromosome or plasmid, a separate *\_fit.csv* file is generated, which contains data at 4096-nucleotide resolution.

133 *skplot.csv* contains all the 4096-nucleotide resolution data as one big file for all processed chromosomes and plasmids. The  
 134 parameters are described in table 1.

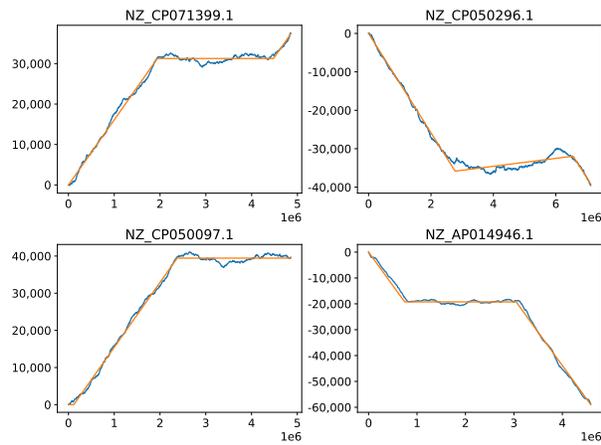
135 *results.csv* meanwhile contains the details of the fits. In this table 2, all marked out squares exist. The actual fields are  
 136 called *alpha1gc*, *alpha2gc*, *gcRMS*, *alpha1ta*, *alpha2ta* etc. DNA sequence shift and div are also specified, and they come from  
 137 the GC skew. *gc0-2*, *ta0-2* refers to codon position. *gcng* and *tang* refer to the non-coding region skews. Finally *sb* denotes the  
 138 strand bias.



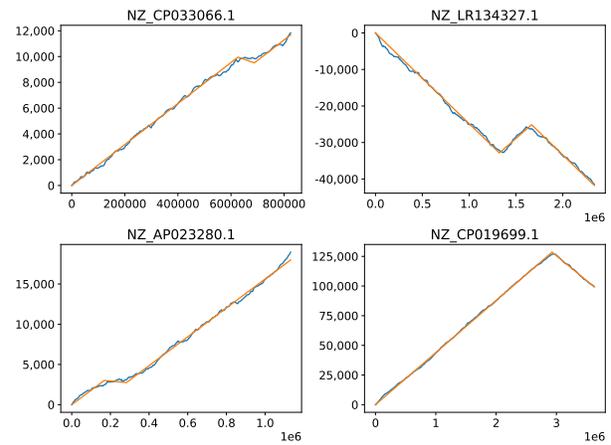
**Figure 3.** Predicted versus actual GC/TA skew for 4093 Firmicutes



**Figure 4.** Scatter graph of codon/strand bias versus GC/TA skew for *C. difficile*



**Figure 5.** Chromosomes with asymmetric skews

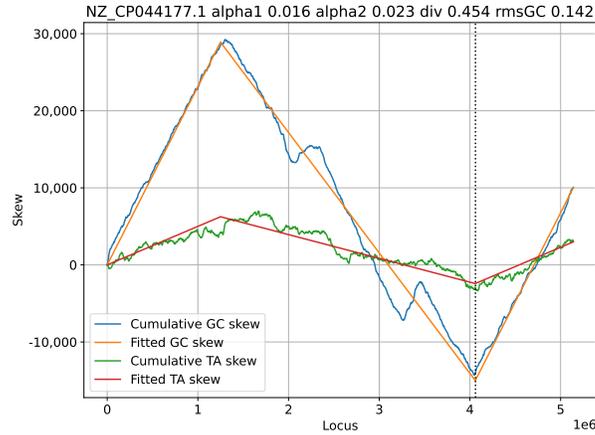


**Figure 6.** Chromosomes with differing strand lengths

	alpha1	alpha2	rms	div	shift
gc	X	X	X	X	X
ta	X	X	X		
gc0	X	X	X		
gc1	X	X	X		
gc2	X	X	X		
ta0	X	X	X		
ta1	X	X	X		
ta2	X	X	X		
gcng	X	X	X		
tang	X	X	X		
sb	X	X	X		

**Table 2.** Skew metrics

139 Table 3 documents the data on codon bias, also split out by leading or lagging strand found in codongc.csv.



**Figure 7.** GC and TA skew for *Salmonella enterica subsp. enterica serovar Concord* strain AR-0407

abspos	locus in chromosome	name	RefSeq ID
accounts0-4	A nucleotide counter	ngcount	Counter of non-coding nucleotides
ccounts0-4	C nucleotide counter	pospos	cumulative positive sense nucleotide counter
gcounts0-4	G nucleotide counter	relpos	relative position within chromosome/plasmid
tcounts0-4	T nucleotide counter	taskew	cumulative TA skew
gcskew	cumulative GC skew	taskew0-3	cumulative TA skew per codon position
gcskew0-3	cumulative GC skew per codon position	taskewNG	cumulative TA skew for non-coding regions
gcskewNG	cumulative GC skew for non-coding regions		

**Table 1.** Fields of skplot.csv

afrac, cfrac, gfrac, tfrac	Fraction of coding nucleotides that are A, C, G or T
leadafrac, leadcfrac, leadgfrac, leadtfrac	Fraction of leading strand coding nucleotides that are A, C, G or T
lagafrac, lagcfrac, laggfrac, lagtfrac	Fraction of lagging strand coding nucleotides that are A, C, G or T
ggcfrac, cgcfrac	The G and C fraction of GC coding nucleotides respectively
atafrac, ttafrac	The A and T fraction of AT coding nucleotides respectively

**Table 3.** Fields in codongc.csv

140 Table 4 documents the fields found in genomes.csv:

fullname	The full chromosome name as found in the FASTA file
acount, ccount, gcount, tcount	Count of A, C, G or T nucleotides
plasmid	Set to 1 in case this sequence is a plasmid
realm1-5	NCBI sourced taxonomic data
protgeneount	Number of protein coding genes processed
stopTAG, TAA, TGA	Number of TAG, TAA and TGA stop codons respectively
stopXXX	Number of anomalous stop codons
startATG, GTG, TTG	Number of ATG, GTG and TTG start codons respectively
startXXX	Number of unusual start codons
dnaApos	position of DnaA gene (not DnaA box!) in the DNA sequence. -1 if not found.

**Table 4.** Fields in genomes.csv

141 Finally, the individual \_fit.csv files contain fields called “Xskew” and “predXskew” to denote the observed X=gc, ta etc  
142 skew, plus the prediction based on the parameters found in results.csv.

## 143 Technical Validation

144 This database models the skews of many chromosomes and plasmids. Validation consists of evaluating the goodness-of-fit  
145 compared to the directly available numbers.

146 The *SkewDB* fits skews to a relatively simple model of only four parameters. This prevents overfitting, and this model has  
147 proven to be robust in practice. Yet, when doing automated analysis of tens of thousands of chromosomes, mistakes will be  
148 made. Also, not all organisms show coherent GC skew.

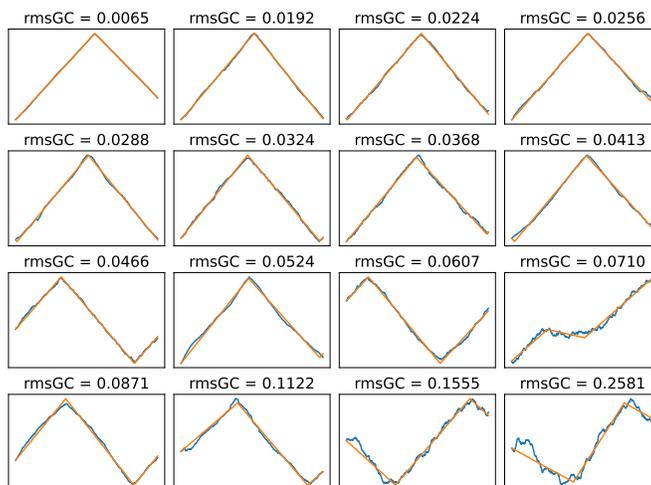


Figure 8. *SkewDB* fits for 16 equal sized quality categories of bacterial chromosomes

149 Figure 8 shows 16 equal sized quality categories, where it is visually clear that the 88% best fits are excellent. It is therefore  
150 reasonable to filter the database on  $RMS_{gc} < 0.16$ . Or conversely, it could be said that above this limit interesting anomalous  
151 chromosomes can be found.

152 The DoriC database<sup>5</sup> contains precise details of the location of the origin of replication. 2267 sequences appear both in  
153 DoriC and in the *SkewDB*. The DoriC origin of replication should roughly be matched by the “shift” metric in the *SkewDB*  
154 (but see Usage notes). For 90% of sequences appearing in both databases, there is less than 5% relative chromosome distance  
155 between these two independent metrics. This is encouraging since these two numbers do not quite measure the same thing.

156 On a similar note, the DnaA gene is typically (but not necessarily) located near the origin of replication. For over 80% of  
157 chromosomes, DnaA is found within 5% of the *SkewDB* “shift” metric. This too is an encouraging independent confirmation of  
158 the accuracy of the data.

159 Finally, during processing numbers are kept of the start and stop codons encountered on all protein coding genes on all  
160 chromosomes and plasmids. These numbers are interesting in themselves (because they correlate with GC content, for example),  
161 but they also match published frequencies, and show limited numbers of non-canonical start codons, and around 0.005%  
162 anomalous stop codons. This too confirms that the analyses are based on correct (annotation) assumptions.

## 163 Usage Notes

164 The existential limitation of any database like the *SkewDB* is that it does not represent the distribution of organisms found in  
165 nature. The sequence and annotation databases are dominated by easily culturable microbes. And even within that selection,  
166 specific (model) organisms are heavily oversampled because of their scientific, economic or medical relevance.

167 Because of this, care should be taken to interpret numbers in a way that takes such over- and undersampling into account.  
168 This leaves enough room however for finding correlations. Some metrics are sampled so heavily that it would be a miracle if  
169 the unculturable organisms were collectively conspiring to skew the statistics away from the average. In addition, the database  
170 is a very suitable way to test or generate hypotheses, or to find anomalous organisms.

171 Finally it should be noted that the *SkewDB* tries to precisely measure the skew parameters, but it makes no effort to pin  
172 down the Origin of replication exactly. For such uses, please refer to the DoriC database<sup>5</sup>. In future work, the *SkewDB* will  
173 attempt to use OriC motifs to improve fitting of this metric.

174 On <https://skewdb.org> an explanatory Jupyter<sup>20</sup> notebook can be found that uses Matplotlib<sup>21</sup> and Pandas<sup>22</sup> to create all the  
175 graphs from this article, and many more. In addition, this notebook reproduces all numerical claims made in this work. The  
176 *SkewDB* website also provides links to informal articles that further explain GC skew, and how it could be used for research.

## Code availability

The *SkewDB* is produced using the Antonie DNA processing software (<https://github.com/berthubert/antonie2>), which is open source. In addition the pipeline is fully automated and reproducible, including the retrieval of sequences, annotations and taxonomic data from the NCBI website. The software has also been deposited with Zenodo<sup>23</sup>.

A GitHub repository is available for this article on <https://github.com/berthubert/skewdb-articles>, which includes this reproducible pipeline, plus a script that regenerates all the graphs and numerical claims from this paper.

## References

1. Frank, A. C. & Lobry, J. R. Asymmetric substitution patterns: a review of possible underlying mutational or selective mechanisms. *Gene* **238**, 65–77 (1999).
2. Marín, A. & Xia, X. GC skew in protein-coding genes between the leading and lagging strands in bacterial genomes: New substitution models incorporating strand bias. *J. Theor. Biol.* **253**, 508–513, <https://doi.org/10.1016/j.jtbi.2008.04.004> (2008).
3. Quan, C.-L. & Gao, F. Quantitative analysis and assessment of base composition asymmetry and gene orientation bias in bacterial genomes. *FEBS Lett.* **593**, 918–925, <https://doi.org/10.1002/1873-3468.13374> (2019).
4. Lu, J. & Salzberg, S. L. SkewIT: The Skew Index Test for large-scale GC Skew analysis of bacterial genomes. *PLOS Comput. Biol.* **16**, e1008439, <https://doi.org/10.1371/journal.pcbi.1008439> (2020).
5. Luo, H. & Gao, F. DoriC 10.0: an updated database of replication origins in prokaryotic genomes including chromosomes and plasmids. *Nucleic Acids Res.* **47**, D74–D77, <https://doi.org/10.1093/nar/gky1014> (2019).
6. O'Donnell, M., Langston, L. & Stillman, B. Principles and concepts of DNA replication in bacteria, archaea, and eukarya. *Cold Spring Harb. Perspectives Biol.* **5**, a010108–a010108, <https://doi.org/10.1101/cshperspect.a010108> (2013).
7. Lilly, J. & Camps, M. Mechanisms of theta plasmid replication. *Microbiol. Spectr.* **3**, <https://doi.org/10.1128/microbiol-spec.plas-0029-2014> (2015).
8. Rudner, R., Karkas, J. D. & Chargaff, E. Separation of *B. subtilis* DNA into complementary strands. 3. Direct analysis. *Proc. Natl. Acad. Sci.* **60**, 921–922, <https://doi.org/10.1073/pnas.60.3.921> (1968).
9. Fariselli, P., Taccioli, C., Pagani, L. & Maritan, A. DNA sequence symmetries from randomness: the origin of the Chargaffs second parity rule. *Briefings Bioinforma.* **bbaa041**, <https://doi.org/10.1093/bib/bbaa041> (2020).
10. Tillier, E. R. & Collins, R. A. The Contributions of Replication Orientation, Gene Direction, and Signal Sequences to Base-Composition Asymmetries in Bacterial Genomes. *J. Mol. Evol.* **50**, 249–257, <https://doi.org/10.1007/s002399910029> (2000).
11. Zhang, R. & Zhang, C.-T. A Brief Review: The Z-curve Theory and its Application in Genome Analysis. *Curr. genomics* **15**, 78–94, <https://doi.org/10.2174/1389202915999140328162433> (2014). Publisher: Bentham Science Publishers.
12. Charneski, C. A., Honti, F., Bryant, J. M., Hurst, L. D. & Feil, E. J. Atypical AT Skew in Firmicute Genomes Results from Selection and Not from Mutation. *PLOS Genet.* **7**, e1002283, <https://doi.org/10.1371/journal.pgen.1002283> (2011).
13. Grigoriev, A. Analyzing genomes with cumulative skew diagrams. *Nucleic Acids Res.* **26**, 2286–2290, <https://doi.org/10.1093/nar/26.10.2286> (1998).
14. Roten, C.-A. H. Comparative genomics (CG): a database dedicated to biometric comparisons of whole genomes. *Nucleic Acids Res.* **30**, 142–144, <https://doi.org/10.1093/nar/30.1.142> (2002).
15. Zhang, C.-T., Zhang, R. & Ou, H.-Y. The z curve database: a graphic representation of genome sequences. *Bioinformatics* **19**, 593–599, <https://doi.org/10.1093/bioinformatics/btg041> (2003).
16. Thomas, J. M., Horspool, D., Brown, G., Tcherepanov, V. & Upton, C. GraphDNA: a java program for graphical display of DNA composition analyses. *BMC Bioinforma.* **8**, <https://doi.org/10.1186/1471-2105-8-21> (2007).
17. Grigoriev, A. Analyzing genomes with cumulative skew diagrams. *Nucleic Acids Res.* **26**, 2286–2290, <https://doi.org/10.1093/nar/26.10.2286> (1998).
18. Nelder, J. A. & Mead, R. A simplex method for function minimization. *The Comput. J.* **7**, 308–313, <https://doi.org/10.1093/comjnl/7.4.308> (1965).
19. Hubert, B. Skewdb: A comprehensive database of gc and 10 other skews for over 28,000 chromosomes and plasmids. *Dryad* <https://doi.org/10.5061/DRYAD.G4F4QRFR6> (2021).

- 224 **20.** Kluyver, T. *et al.* Jupyter notebooks – a publishing format for reproducible computational workflows. In Loizides, F. &  
225 Schmidt, B. (eds.) *Positioning and Power in Academic Publishing: Players, Agents and Agendas*, 87 – 90 (IOS Press,  
226 2016).
- 227 **21.** Hunter, J. D. Matplotlib: A 2d graphics environment. *Comput. Sci. & Eng.* **9**, 90–95, <https://doi.org/10.1109/MCSE.2007.55>  
228 (2007).
- 229 **22.** Reback, J. *et al.* pandas-dev/pandas: Pandas 1.3.2. *zenodo* <https://doi.org/10.5281/zenodo.5203279> (2021).
- 230 **23.** Hubert, B. & Beaumont Lab. berthubert/antonie2: Skewversion 1.0. *zenodo* <https://doi.org/10.5281/ZENODO.5516524>  
231 (2021).
- 232 **24.** Hol, F. J. H., Hubert, B., Dekker, C. & Keymer, J. E. Density-dependent adaptive resistance allows swimming bacteria to  
233 colonize an antibiotic gradient. *The ISME J.* **10**, 30–38, <https://doi.org/10.1038/ismej.2015.107> (2016).

## 234 **Acknowledgements**

235 I would like to thank Bertus Beaumont for helping me to think like a biologist, and Jason Piper for regularly pointing me to the  
236 relevant literature. In addition, I am grateful that Felix Hol kindly allowed me to field test my software on his DNA sequences<sup>24</sup>.  
237 Twitter users @halvorz and @Suddenly\_a\_goat also provided valuable feedback.

## 238 **Author contributions statement**

239 B.H. did all the work.

## 240 **Competing interests**

241 The author declares no competing interests.

## 242 **Figures & Tables**

243 Figures:

- 244 1. Sample graph showing *SkewDB* data for *Lactiplantibacillus plantarum* strain LZ95 chromosome
- 245 2. Scatter graph of 25,000 chromosomes by superphylum, GC skew versus TA skew
- 246 3. Predicted versus actual GC/TA skew for 4093 Firmicutes
- 247 4. Scatter graph of codon/strand bias versus GC/TA skew for *C. difficile*
- 248 5. Chromosomes with asymmetric skews
- 249 6. Chromosomes with differing strand lengths
- 250 7. GC and TA skew for *Salmonella enterica subsp. enterica serovar Concord* strain AR-0407
- 251 8. *SkewDB* fits for 16 equal sized quality categories of bacterial chromosomes

252 Tables:

- 253 1. Fields of skplot.csv
- 254 2. Skew metrics
- 255 3. Fields in codongc.csv
- 256 4. Fields in genomes.csv