

1 Linguistic laws in biology

2

3 Stuart Semple^{1*}, Ramon Ferrer-i-Cancho^{2*} and Morgan L. Gustison^{3*}

4

5 ¹ School of Life and Health Sciences, University of Roehampton, London, UK

6 ² Complexity and Quantitative Linguistics Laboratory, Laboratory for Relational Algorithmics, Complexity, and
7 Learning Research Group, Departament de Ciències de la Computació, Universitat Politècnica de
8 Catalunya, 08034 Barcelona, Catalonia, Spain

9 ³ Department of Integrative Biology, University of Texas at Austin, Austin, TX, USA

10 * All authors contributed equally

11

12 **Corresponding author:** Semple, S. (s.semple@roehampton.ac.uk)

13

14 **Key words:** Zipf, Menzerath, communication, compression, information theory

15

16 **Abstract**

17 Linguistic laws, the common statistical patterns of human language, have been investigated by
18 quantitative linguists for nearly a century. Recently, biologists from a range of disciplines have
19 started to explore the prevalence of these laws beyond language, finding patterns consistent with
20 linguistic laws across multiple levels of biological organisation - from molecular (genomes, genes
21 and proteins) to organismal (animal behaviour) to ecological (populations and ecosystems). We
22 propose a new conceptual framework for the study of linguistic laws in biology, comprising and
23 integrating distinct levels of analysis - from description to prediction to theory building. Adopting
24 this framework will provide critical new insights into the fundamental rules of organization
25 underpinning natural systems, unifying linguistic laws and core theory in biology.

26 **The laws of language – from linguistics into biology**

27 Investigations of linguistic laws, the common statistical patterns of human language, have been a
28 cornerstone of quantitative linguistics since the first such laws were proposed in the 1930's [1–3].
29 More recently, these laws have attracted attention from biologists, who are beginning to explore
30 their prevalence beyond human language (Table 1, Figure 1). Work in this area initially focused on
31 Zipf's rank-frequency law (more often known as Zipf's law), which states that the relative frequency
32 of a word is inversely proportional to its frequency rank [4]. In recent years, there has been a
33 notable increase in research on two further linguistic laws - Zipf's law of abbreviation, which
34 describes the tendency of more frequently used words to be shorter [4], and Menzerath's law,
35 which states that the larger the linguistic construct, the smaller its constituent parts [2,5,6]. A wide
36 range of further linguistic laws exists, and many are yet to be investigated beyond language (e.g.
37 [7]).

38
39 Biologists have tested for patterns consistent with linguistic laws among diverse taxa and different
40 levels of biological organisation - molecular [8–17], organismal [18–27] and ecological [28–31] (see
41 Table 1). This work necessitates more general formulations of these laws, which originally were
42 framed in specific linguistic constructs. Additionally, researchers need to reconsider not just the
43 appropriate comparable elements of analysis (e.g. call types, instead of words) but also the
44 relevant dimensions to be measured (e.g. call type duration, instead of word length in letters). The
45 expansion of studies of linguistic laws beyond language raises the question of whether we should
46 rename these phenomena - e.g. to biological laws - to reflect their broader applicability. However,
47 we feel it is appropriate for now to retain the original terminology, in recognition of the field in which
48 the laws were first derived.

49
50 Here, we advocate a new approach to exploring linguistic laws in biology, which moves
51 investigations from the level of description to the levels of explanation, prediction and theory
52 building.

53

54 **Studies of linguistic laws in biological systems**

55 Zipf's rank-frequency law

56 The relationship between the frequency of a word and its rank (1st most frequent, 2nd most frequent
57 etc.) that is expected by Zipf's rank-frequency law [4] (Figure 1B), is traditionally explored in written
58 texts and spoken languages [32,33]. Rank and frequency are by definition negatively associated,
59 and linguists studying this law (or the closely related modification, Zipf-Mandelbrot law – see Table
60 1) focus on characterizing the degree of linearity and steepness (exponent) of the log-log slope
61 between these measures. In written English, the relationship is highly linear with an exponent ~1
62 [34], and a study of this law across 100 typologically diverse languages found a narrow exponent

63 range centred around 1 but with some variation (0.76-1.44) [35]. Studies of Zipf's rank-frequency
64 law beyond language come from different levels of biological organisation and diverse taxa (Figure
65 1C, Table 1).

66
67 At the molecular level, frequency distributions consistent with this law (or its sister law, Zipf's
68 number-frequency law – Table 1), are documented in systems including oligonucleotides in DNA
69 sequences [9], secondary structures of RNA [10], gene expression [36] and the size of protein
70 superfamilies [37]. At the organismal level, most work focuses on animal communication, with
71 patterns predicted by this law (or Mandelbrot's modification) documented in the vocalisations of
72 birds (e.g. [21,38–40]) and of both terrestrial and marine mammals [41–43], as well as in primate
73 gestural communication [20]. Beyond the organismal level, Zipf's rank-frequency law describes
74 abundance distributions of both species and communities. For example, patterns consistent with
75 this law are found in the spatial distribution of land snails, *Vallonia pulchella* [28], and the relative
76 abundance of species within a broad range of plant and animal communities [44]. This law has
77 also been applied to epidemiological data, for example to explore patterns of geographic
78 distribution of COVID-19 [30]. In sum, Zipf's rank-frequency law is a notable example of how
79 linguistics laws can apply to biological systems beyond language.

80

81 Zipf's law of abbreviation

82 A significant negative relationship between word length and frequency of use, as expected under
83 Zipf's law of abbreviation [1] (Figure 1B, Table 1), is found across hundreds of typologically diverse
84 languages in written form [45], as well as in spoken [46–48] and sign language [49]. Beyond
85 language, and in marked contrast to Zipf's rank-frequency law, studies that explicitly test Zipf's law
86 of abbreviation have been conducted only at the organismal level - in the behaviour (and
87 particularly vocal communication) of birds and mammals - with the number of studies increasing
88 noticeably in recent years (Figure 1C). A negative relationship with frequency of occurrence has
89 been found for the length of song bouts (measured as the number of calls) in black-capped
90 chickadees, *Poecile atricapillus* [50], call type duration in Formosan macaques, *Macaca cyclopis*
91 [51], vocal phrase size (in component units) in Indri, *Indri indri* [52], gesture duration in
92 chimpanzees, *Pan troglodytes* [23] and the size (in component units) of surface behavioural
93 patterns of dolphins, *Tursiops truncatus* [24]. However, not all species tested follow Zipf's law of
94 abbreviation, at least in the initial analyses [23,53,54]; these examples provide important insights
95 into where and how this law should be explored.

96

97 It is important to note that the absence of studies explicitly exploring Zipf's law of abbreviation at
98 molecular and ecological levels does not mean that the patterns predicted by the law have not
99 been investigated at those levels. For example, studies of global size-density relationships in
100 ecology have documented a negative association between population density and body size of

101 species [55] - more frequently occurring species tend to be smaller. This finding highlights a key
102 point: that linguistic laws and established concepts in different fields of biology may describe
103 similar, or even identical, patterns. Identifying such common ground provides a foundation for
104 unification across disciplines, and for the development of more general theory about the
105 organisation of biological systems.

106

107 Menzerath's law

108 Patterns consistent with Menzerath's law - a negative relationship between the size of the whole
109 and the size of the constituent parts [5] (Figure 1B, Table 1) - occur across many languages and
110 linguistic levels (e.g. phonemes, syllables, words, clauses and sentences [56,57]). This law is
111 documented not just in written texts, but also in spoken [47,48] and signed language [58].
112 However, some studies fail to find the patterns predicted by Menzerath's law [59,60], suggesting it
113 is not as widespread in human language as Zipf's law of abbreviation. Such inconsistencies may
114 reflect variation in how the 'whole' and 'constituent parts' are connected; it is thought that these
115 linguistic levels should be immediately consecutive (e.g. clauses and phrases), rather than distant
116 (e.g. clauses and words) [61]. Investigations of Menzerath's law (or the closely related Menzerath-
117 Altmann law – Table 1) in biological systems initially focused on genes, genomes and proteins, but
118 over the last five years there has been growing interest in applying this law to different forms of
119 animal communication (Figure 1C). At the molecular level, studies demonstrate a negative
120 relationship between exon number and size in genes [14,62], domain number and size in proteins
121 [16], segment number and size in RNA [17], and chromosome number and size in genomes [8,15].
122 In animal communication, patterns consistent with Menzerath's law are seen in primate vocal
123 sequences [26,52,63–65] and bird song [22,27], as well as in the gestural communication of
124 chimpanzees [23]. However, vocal sequences of a number of primate and bird species do not
125 follow this law [27,66].

126

127 Although there is an absence of explicit tests of Menzerath's law at the ecological level, patterns
128 consistent with this law have in fact been documented there. For example, studies of cross-
129 community scaling relationships show a negative association between the total number of
130 individuals in an ecological community and the average size of an individual in that community [55]
131 – the larger the whole, the smaller the constituent parts. This result again highlights the existence
132 of common ground between linguistic laws and - on the surface, completely unrelated - core
133 concepts in other fields.

134

135 Other linguistic laws

136 In addition to work exploring these three linguistic laws in biological systems, a small number of
137 studies investigated other linguistic laws beyond language (Table 1). Patterns consistent with
138 Herdan's law, which describes the relationship between the number of unique words in a text and

139 the overall size of the text, exist in the proteomes of viruses, Archaea, bacteria and Eukarya [13].
140 Zipf's meaning-frequency law, which states that more frequently occurring words tend to have
141 more meanings, is found in chimpanzee gestures [25] and bottlenose dolphin whistles [19].
142 Additional tests of these two linguistic laws, and others that have yet to be explicitly tested in
143 biological systems (see [7] for examples), offer opportunities to explore commonalities between
144 human language and diverse natural systems. As described above, it is also possible that patterns
145 predicted by linguistic laws have already been investigated in biological systems but without
146 explicit reference to (and perhaps without knowledge of) one of these laws.

147 148 Conditions for the investigation of linguistic laws in biological systems

149 Linguistic laws are inherently simple in design (Figure 1A and B). Those described above share
150 one core criterion: that a data structure contains at least one group of discrete units. This is the
151 only criterion for Zipf's rank-frequency law and Herdan's law, which explore the distribution
152 patterns of categorically distinct units [4,7,67]. Zipf's law of abbreviation requires, in addition, that
153 units are measured by some aspect of their size, while Zipf's meaning-frequency law requires that
154 units have one or more distinct meanings (or proxies thereof). Menzerath's law involves a
155 hierarchical data structure of two or more unit levels, where larger constructs are made up of units
156 measured by size. All these criteria are broadly applicable across biological systems, which are
157 generally composed of discrete parts and organized in a hierarchical structure.

158
159 Beyond the satisfaction of basic criteria, a key issue remains in the exploration of certain linguistic
160 laws in biological systems: what are the most appropriate 'currencies' to use in place of the specific
161 linguistic measures from the original formulation of these laws? In studies of Zipf's law of
162 abbreviation and Menzerath's law in animal communication, for example, duration of call or gesture
163 types is often assessed, but this approach most likely reflects that duration is the easiest - rather
164 than the most appropriate - measure of unit size. A study of rock hyrax, *Procavia capensis*,
165 vocalisations indicates that the choice of currency can be critical [68]. Among males, a negative
166 relationship was found between call type duration and frequency of use, in line with Zipf's law of
167 abbreviation, but among females this relationship was positive (i.e. opposite to the law). However,
168 when call amplitude rather than duration was explored, both sexes conformed to the law - more
169 frequently used call types were quieter. These results highlight that simple, discrete measures of
170 size may not always be appropriate or sufficient. Instead, broader and perhaps multidimensional
171 measures of magnitude, or more precise indices of energetic cost, may be more suitable.

172
173 The need to identify a suitable currency may limit the range of systems in which particular linguistic
174 laws can be explored. For example, laws of word meaning can be tested in those animal
175 communication systems where signal types can be assigned a specific meaning (e.g. chimpanzee
176 gestures [25]), and perhaps also in genetic systems where, for example, different combinations of

177 codons ‘mean’ (i.e. translate into) different amino acids. However, there seems no obvious analogy
178 for meaning at an ecological level, so extending such laws to populations and ecosystems may not
179 be feasible. Ultimately, the appropriate currency to use when testing linguistic laws in biology
180 should be informed by a strong underpinning theory, from which clear predictions are derived to
181 describe the expected patterns and how these patterns are to be assessed.

182

183 The universality of linguistic laws, and the role of ‘exceptions’

184 A key issue around the investigation of linguistic laws, in human language and in other systems,
185 relates to the notion of universality. In traditional linguistics [69], two types of ‘universals’ are
186 typically considered. Absolute universals describe exception-less patterns, while statistical
187 universals describe patterns occurring significantly above chance levels. The latter are the
188 universals that should be considered in relation to linguistic laws, but in studies of these laws
189 beyond human language the distinction between these two types of universals is perhaps not so
190 readily recognised. Finding patterns that do not fit those predicted by linguistic laws - ‘exceptions’ -
191 is a common theme (e.g. [54,66,70]). However, such exceptions are certainly not unexpected for
192 statistical universals, and can in fact be important for highlighting alternative ways or parts of the
193 system where laws can be tested.

194

195 For example, a study of Zipf’s law of abbreviation in four bat species found a negative correlation
196 between duration and frequency of production for short range social calls, but not for distress calls
197 [71]. Similarly, while Bezerra et al [54] found no link between duration and frequency of use of call
198 types in the full vocal repertoire of common marmosets, *Callithrix jacchus*, a subsequent analysis
199 by Ferrer-i-Cancho and Hernández-Fernández [53] found conformity to the law in a subset of the
200 repertoire comprising close range social calls. These studies indicate that investigations of
201 linguistic laws should consider patterns that might occur in subsets of the whole repertoire.
202 Developing stronger theory around linguistic laws will allow for specific predictions to be made
203 about where (i.e. in which subsets) we might expect to find conformity to these laws.

204

205 A study of Menzerath’s law in gorilla, *Gorilla gorilla*, vocal sequences [72] highlights another
206 consideration - whether we should include situations where a constituent part represents the whole
207 construct. In this study, the predicted negative relationship between call duration and sequence
208 length (in terms of the number of constituent calls) was found, but only when analyses included
209 sequences consisting of a single call. When analyses considered only sequences of two or more
210 calls, the negative relationship no longer remained, and in fact a weak positive association was
211 found. Future studies of Menzerath’s law should consider data both with and without sequences of
212 length one, to explore this phenomenon further (notably, studies of Menzerath’s law in language
213 often exclude words of one syllable [48,60]).

214

215 **Studying linguistic laws in biology: a new framework**

216 As investigations of linguistic laws in molecular biology, organismal biology and ecology increase in
217 number (Figure 1C), it is important now to ask what these laws mean in different systems, why
218 these patterns occur, and how we can move the field forward to provide the most important
219 insights. To this end, we propose a new framework for the study of linguistic laws in biology (Figure
220 2) that builds on the foundations and aims of the scientific method [73,74]; this framework is
221 underpinned by distinct analytical levels and the integration of research across levels. To date,
222 most studies fit into the lower levels of this framework, involving exploratory and descriptive
223 statistics (Level 1) and mathematical modelling (Level 2) or (less often) inferential and predictive
224 modelling (Level 3). There have also been insightful attempts at theory construction for specific
225 laws or study systems, but extension to the creation of more general theory (Level 4) is typically
226 lacking. An important goal for future research will be to shift from the level-centric approach
227 adopted thus far, to bottom-up and top-down integrative approaches where empirical research
228 informs, and is guided by, generalized theory and explicit hypotheses.

229

230 Level 1. Exploratory and descriptive statistics: The most common approach to investigate linguistic
231 laws in biological data involves the description and exploration of the patterns these laws predict.
232 At the simplest level, researchers have described qualitatively how well the observed data align
233 with the patterns predicted by a specific linguistic law (e.g. [38]). Stronger, quantitative approaches
234 (Box 1) fit curves or distributions to data, or test for predicted correlations between variables
235 (examples are found in Table 1). Description and exploration are critical for understanding when
236 and where patterns consistent with linguistic laws appear, and identifying exceptions to laws allows
237 these to be explored at higher levels in the framework (see Level 3). However, this approach on its
238 own is unable to answer basic questions about how and why laws emerge in the first place. Such
239 questions are well-suited to mathematical modelling approaches that describe and explain law-like
240 patterns.

241

242 Level 2. Descriptive mathematical modelling: This approach is oriented towards describing simple
243 processes that reproduce law-like patterns. A common focus of linguistic law research at this level
244 involves stochasticity. This stems from early observations that Zipf-like word frequency
245 distributions arise from random processes that combine subunits to form units (i.e. Miller's random-
246 monkey model) and generate new text (i.e. Simon's model) [60,75–78], and that other stochastic
247 processes (e.g. random breakage) also reproduce patterns that look like Menzerath's law [79].
248 Such studies call into question the meaningfulness of linguistic laws, but these arguments are
249 refuted by multiple lines of evidence (see Box 2). A complementary line of work tests whether the
250 expression of law-like patterns by these stochastic models accurately captures patterns in real-
251 word data [60,80–82]. On the whole, descriptive modelling is useful for sorting out how law-like

252 patterns can be reproduced by stochastic or deterministic mechanisms. This approach is limited,
253 however, in that it is unable to answer questions about why such patterns emerge.

254

255 Level 3. Inferential and predictive mathematical modelling: Mathematical modelling is not always
256 descriptive, but critically can also serve to interpret, explain, and generate predictions (see Box 2).
257 A noteworthy example comes from work combining mathematical models from information theory
258 and language evolution to explain how arbitrary signals become associated with meaning [83].
259 Mathematical work also links information theory to linguistic laws. Ferrer-i-Cancho et al [84,85]
260 demonstrated a mathematical relationship between Zipf's law of abbreviation and the information
261 theoretic principle of compression. This is the principle of minimizing the average length of an
262 element in a system, for example by inversely aligning the length of elements with their frequency
263 of occurrence. Gustison et al [26] and Ferrer-i-Cancho et al [86] similarly used mathematical
264 models to connect Menzerath's law and Zipf's rank-frequency law to compression. This work,
265 moving beyond the level of description and simple generative processes, strongly suggests that
266 linguistic law-like patterns in certain systems reflect selection for coding efficiency; from here,
267 predictions under alternative hypotheses can be developed and tested.

268

269 For example, in animal communication there is a trade-off between signalling efficiency (linked to
270 compression) and the need to accurately transmit information [84]. This leads to the prediction that
271 conformity with these linguistic laws is expected to be strong in short range communication, but
272 less prevalent in long-range signalling where redundancy helps ensure that information is
273 accurately received [26,84]. This prediction has been supported by studies finding that conformity
274 to laws is not seen in analyses (Level 1) of complete repertoires, but is seen in analyses of close
275 range signals [53,71]. Mathematical approaches that generate predictions in this way are powerful
276 (the stochastic generative processes in Level 2 lack this crucial property – see Box 2) and provide
277 the groundwork for building general theory.

278

279 Level 4. General theory: Biological data are rife with patterns, but patterns are of limited value
280 without a solid theoretical framework to explain them [87]. An ideal framework, i.e. a scientific
281 theory, is constructed from a set of interconnected principles that are independent from the
282 observable phenomena [74]. With the growing body of exploratory, descriptive, and inferential
283 research on linguistic laws in biology, alongside the recognition that core concepts in certain fields
284 make strikingly similar predictions to those of linguistic laws, we can now begin to build fully-
285 fledged scientific theory. Common themes emerging from different studies provide a key starting
286 point for this endeavour. For example, studies of animal communication and of gene structure
287 suggest that patterns consistent with linguistic laws in these systems reflect evolutionary trade-offs
288 - between the costs and benefits of signal compression in animal communication [26], and
289 between the costs and benefits of structural change in genes [14]. A general theory about the

290 organization of biological systems, that is underpinned by evolutionary selection and that can
291 accommodate these and other trade-offs, would thus be applicable across molecular and
292 organismal levels of biological organisation.

293

294 We have developed one such theory, building from an empirical and mathematical exploration of
295 the vocal sequences of geladas, *Theropithecus gelada* [26]. Here, a negative correlation was found
296 between vocal sequence size and the duration of constituent calls, in line with Menzerath's law
297 (Level 1). Sequence lengths were described mathematically as being consistent with a
298 'memoryless' process (Level 2), and Menzerath's law was interpreted as a prediction of
299 compression via a generalized cost function (Level 3). These findings were combined with previous
300 work on this law in genes, genomes and proteins, to develop the theory (Level 4) that compression
301 - reflecting a trade-off between efficiency of coding and information transmission success -
302 underpins not just animal (including human) communication but also biological information systems
303 in the broadest sense. Such bottom-up development of general theory then facilitates a top-down
304 approach, driving empirical work and hypothesis development at lower levels. For example, James
305 et al [27] built on this general theory to explore how production mechanisms and learning
306 contribute to the emergence of Menzerath's law in bird song, and proposed the hypothesis that
307 ease of motor production underpins compressional organization in this communication system.

308

309 An important challenge now is to develop theory that can incorporate exploration of linguistic laws
310 across all biological levels, and thus can unite evolutionary and ecological studies. A vital first step
311 in developing such theory is to reframe linguistic laws in a more abstract way. This abstraction
312 relies on identifying appropriate currencies that are applicable across a wide range of biological
313 systems. Here again, identifying common ground from studies in diverse systems can provide a
314 valuable starting point. For example, the finding that species abundance distributions follow Zipf's
315 rank-frequency law may reflect the differential 'cost of species' [44], where cost reflects the amount
316 of energy required to produce and maintain the organism (e.g. carnivores are 'costly' due to their
317 high trophic level and thus rare). In human and animal communication, patterns consistent with
318 Zipf's rank-frequency law may similarly reflect a differential cost [88] - here not of species but of the
319 various components of the repertoire (for example, in terms of the energetic costs of production or
320 perception).

321

322 We propose reframing linguistic laws in terms of a generalised cost function based on energy, and
323 tentatively propose the start of an overarching general theory: that patterns consistent with
324 linguistic laws reflect pressures that shape the allocation of finite energy in a system. The
325 pressures involved will differ markedly between systems, and identifying and exploring the exact
326 underpinning mechanisms will be an important step. At the molecular and organismal levels,
327 evolutionary selection for optimal allocation of constrained energy may be key; then at the

328 ecological level, competition within and between species for the finite available energy in the
329 environment may be critical. Further development and refinement of a general theory of this kind,
330 based on energy and its allocation, is a promising approach to unify all levels of biological
331 organisation. Recent work in molecular biology [89], organismal biology and ecology [90] has
332 similarly focused on energy as a unifying currency, and an exciting prospect is to reconcile general
333 theory around linguistic laws with, for example, metabolic theory of ecology [91] and the equal
334 fitness paradigm [92].

335

336 Level integration: Contemporary studies of linguistic laws in biology are largely level-centric, with
337 an emphasis on the descriptive approaches outlined in Levels 1 and 2 above. These lower-level
338 studies are important for developing basic knowledge, but their narrow focus means that we risk
339 being left with a patchwork of datasets and models for different laws and biological systems. To
340 advance understanding, we need more research that connects and integrates levels, achieving a
341 cyclical process (akin to that of the scientific method [74]) in which bottom-up and top-down
342 approaches work in tandem to build, test, and refine general theory (Figure 2). This process should
343 involve research on systems that conform to, or are an exception to (e.g. [54,66,70]), one or more
344 linguistic laws. Exceptions can be used to make clear predictions for testing in independent
345 systems, and in turn, drive theory building and refinement. Collectively, integrative research efforts
346 will promote general theory that joins distinct, unifying principles to understand how, when, and
347 why linguistic laws manifest in biological systems. Such theory will be particularly powerful,
348 allowing the development of clear hypotheses and predictions that are sufficiently abstract in their
349 framing (for example in terms of the system, or units of analysis) as to be applicable across all
350 levels of biological organisation.

351

352 **Concluding remarks**

353 Studies of linguistic laws in biology have provided important insights, but the true power of such
354 work has yet to be fully recognised. Many questions remain to be addressed (see Outstanding
355 Questions), and there are opportunities to expand this line of inquiry into new disciplines and
356 towards linguistics laws that have yet to be investigated beyond human language. There is
357 potential to analyse existing biological databases to test linguistic laws, and to explore common
358 ground between such laws and core principles and concepts in different scientific fields. This
359 opens up exciting opportunities for broad scale comparative analyses – not just across taxa but
360 also across diverse biological systems and different levels of biological organisation. Above all,
361 there is an opportunity to reshape how we investigate linguistic laws in biology, integrating different
362 levels of analysis with the ultimate goal of generating and testing general theory. Such unifying
363 work, built on the open exchange and cross-fertilisation of ideas from multiple disciplines, will
364 provide new understanding of the fundamental rules of organization underpinning diverse
365 biological systems – from molecular to organismal to ecological.

366

367 **Acknowledgments**

368 We sincerely thank Daniel Perkins, Logan James, Daniel Takahashi, Antoni Hernandez-
369 Fernandez, two anonymous reviewers and Andrea Stephens for their insightful comments. MLG is
370 supported by the grant K99MH126164 from NIH (National Institutes of Health). RFC is supported
371 by the grant TIN2017-89244-R from MINECO (Ministerio de Economía, Industria y Competitividad)
372 and by the recognition 2017SGR-856 (MACDA) from AGAUR (Generalitat de Catalunya).

373 **Box 1. Key methodological issues in the study of linguistic laws**

374 Two main statistical approaches are used to investigate linguistic laws: model fitting and
375 correlation. Model fitting, which is the traditional approach taken in quantitative linguistics, involves
376 finding a mathematical function (i.e. an equation) that defines the dependency between two
377 variables. Putative mathematical functions for a number of linguistic laws are shown in Table 1. A
378 key advantage of the model fitting approach is that it helps characterize the non-linear behaviour of
379 a law. Any resulting equations are then tailored by modifying parameters to fit new datasets (e.g.
380 Zipf's rank-frequency exponent varies across languages [35]). On the other hand, these fitted
381 equations are only first approximations to the true, highly complex, relationship between variables
382 [34,93,94]. The validity of many of these equations has been questioned using both theoretical and
383 empirical arguments [32,95–97]. As a reaction to these challenges, as well as a need to generalize
384 laws to diverse systems [84], correlational approaches have become increasingly popular
385 (e.g. [45]).

386

387 A correlational approach involves the use of statistical tools to test for positive or negative
388 associations between variables. This approach has some key advantages. First, it simplifies the
389 analysis with respect to curve fitting, and the use of appropriate statistical tools (e.g. Spearman's
390 correlation, linear mixed effect models) removes some assumptions about how variables are
391 structured and related to one another, or controls for the effect of multiple factors. Second, it is
392 supported by theoretical arguments involving optimal coding that predict a negative (or null)
393 correlation consistent with Zipf's law of abbreviation or Menzerath's law [85]. A limitation of
394 correlational analysis is that it is not appropriate for all laws, specifically those where variable
395 definitions are intrinsically related (e.g. frequency can only decrease as rank increases in Zipf's
396 rank-frequency law; type number can only increase with token number in Herdan's law; Table 1).
397 For other laws, it is important to confirm that there are not 'inevitable' correlations due to variables
398 being functionally dependent [79,98]. Researchers have used several strategies to successfully
399 address this criticism for Menzerath's law and Zipf's law of abbreviation [23,26,99].

400

401 There are many analytical tools available to explore linguistic laws while avoiding statistical pitfalls.
402 In addition to the references above, we also point readers to discussions on the advantage of
403 maximum likelihood over least squares [100], the use of randomization approaches to develop
404 appropriate null models [27], the problem of rank as a random variable [97], and the problem of
405 independence violation [93].

406 **Box 2. Debates on the meaningfulness of linguistic laws**

407 Since Zipf's foundational research [4], many researchers have cast doubts on the meaningfulness
408 and utility of linguistic laws. A recurrent criticism is that linguistic laws are inevitable. For Zipf's
409 rank-frequency law, this argument is illustrated with the metaphor of typing, whereby choosing
410 simple units (e.g. letters, nucleotides) at random can result in a data structure conforming to the
411 law [75,78,101,102]. For Menzerath's law, this argument is derived from the definition of the size of
412 the constituent parts as an average [79]: if the total size of a whole construct is constrained, then
413 randomly breaking it into a few versus many pieces will inevitably result in larger versus smaller
414 parts. This line of reasoning also extends to Zipf's law of abbreviation [99]. However, the
415 inevitability of such patterns has been falsified in several ways: (i) by finding patterns across
416 different species and systems that are not consistent with laws (e.g. [23,43,103]); (ii) through
417 mathematical analysis to show that random typing does not in fact reproduce Zipfian laws of unit
418 frequency [104] and that defining constituent part size as an average does not inevitably lead to
419 Menzerath's law [98]; (iii) by showing statistical differences in how laws are expressed in real data
420 as opposed to the artificially constructed datasets used to support the inevitability argument
421 [80,81,105].

422
423 Another criticism, applied mainly to Zipf's rank-frequency law but that can be generalized to other
424 laws, is that the presence of a law does not allow inferences to be made about how it works or
425 what function it serves in a specific system [60,106,107]. The main reason for this criticism is that
426 Zipf's rank-frequency law (and others) can be reproduced in many ways, which questions whether
427 a law is meaningful for any specific system. This criticism is an important one, and reflects the
428 current emphasis on exploratory and descriptive studies (Levels 1 and 2 in Figure 2). To
429 thoroughly address this criticism, a shift is needed towards work that makes explanations for laws,
430 develops predictions that can be tested in independent contexts, and builds a general theoretical
431 framework to guide how future studies are designed (Levels 3 and 4). An example of a novel
432 prediction is provided by a model originally designed to reproduce Zipf's rank-frequency law based
433 on cognitively realistic principles; this model sheds light on how children learn new words, while
434 random typing or Simon's model fail to make any prediction [82]. The multiplicity of ways of
435 reproducing a law (Level 2), e.g. [87], does not imply a multiplicity of explanations (Levels 3 and 4).

436 **Figure legends**

437

438 **Figure 1. Common linguistic laws and publishing trends.** (A) An artificial dataset is illustrated
439 using different symbol-colour combinations. There is a repertoire of eight discrete unit types that
440 differ in their relative sizes (left), and the full dataset involves a collection of 64 units that are
441 grouped into aggregates ranging in size from one to ten units (right). For example, in human
442 language the units could be words and the aggregates phrases, while in animal communication the
443 units could be call types in the vocal repertoire and the aggregates vocal sequences. (B) This
444 artificial dataset conforms to patterns expected under Zipf's rank-frequency law (left), Zipf's law of
445 abbreviation (centre), and Menzerath's law (right). The symbols in the Zipfian law plots represent
446 specific unit types, while symbols in the Menzerath's law plot represent aggregates. (C) Cumulative
447 number of papers that explicitly test these three linguistic laws in biological systems, published
448 over the last ten years, assessed from Web of Knowledge using search terms 'Zipf*' or
449 'Menzerath*' on February 22, 2021. For Zipf's rank-frequency law, also included are the closely
450 related Zipf-Mandelbrot law and Zipf's number-frequency law; for Menzerath's law, also included is
451 Menzerath-Altmann law. Papers are categorized based on the biological level of organization of
452 the study system (molecular, organismal, ecological). Striking differences are clear between the
453 laws in terms of the total number of studies conducted, and in particular the levels of organisation
454 at which each law has been investigated.

455

456

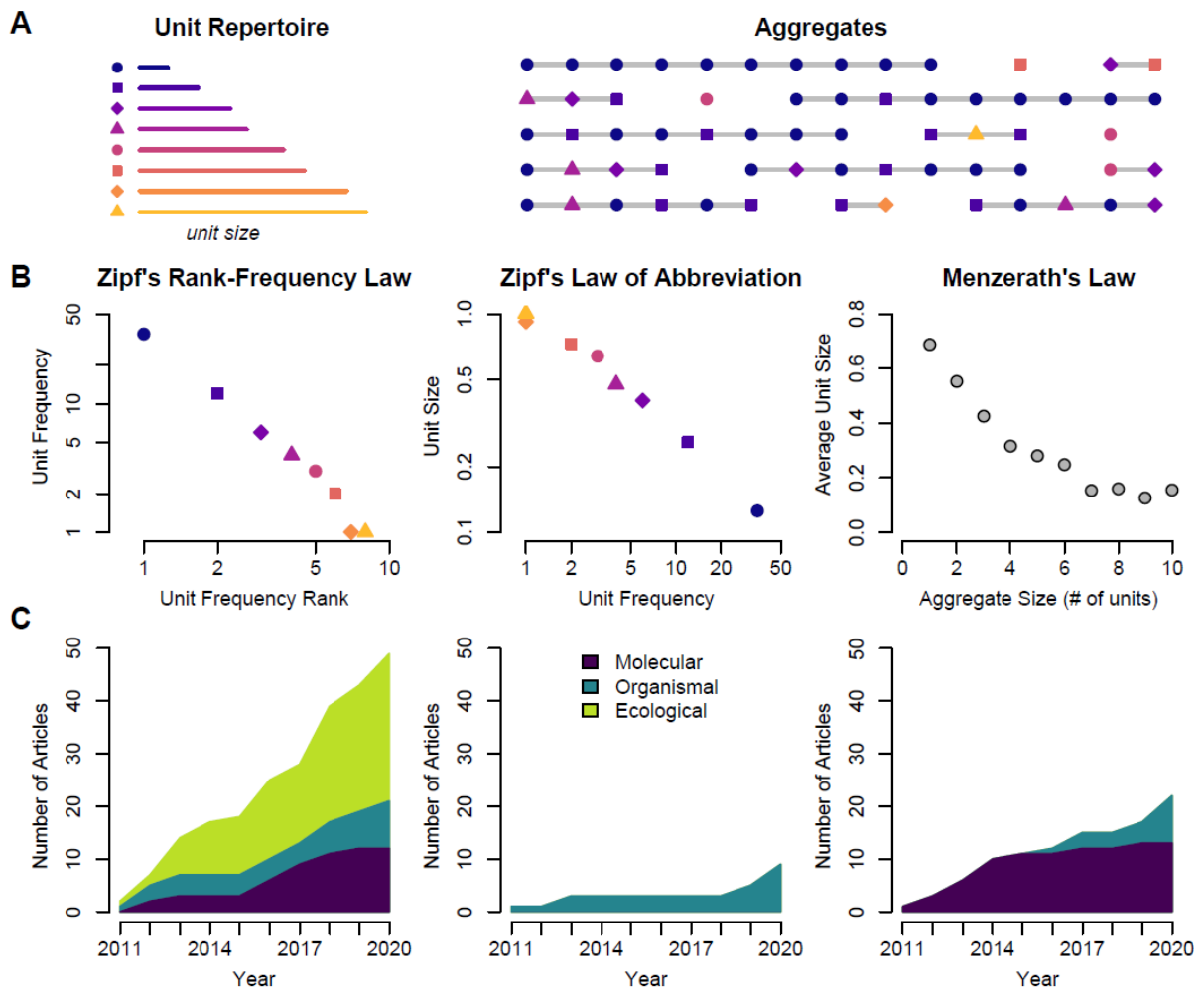
457 **Figure 2. Conceptual framework for investigating linguistic laws in biological systems.**

458 Research on linguistic laws typically falls into one of four analytical levels. 'Exploratory and
459 Descriptive Statistics' (Level 1) is the most basic level, involving studies that test for conformity to
460 linguistic laws in real world biological systems (the signs of correlations and equations in Table 1).
461 'Descriptive Mathematical Modelling' (Level 2) involves work that describes processes that
462 reproduce these laws. 'Inferential and Predictive Modelling' (Level 3) involves research that uses
463 computational approaches to explain possible functions of these laws and develops testable
464 predictions. 'General Theory' (Level 4) involves work to build scientific theory by integrating a set of
465 interconnected principles that are independent from the observable phenomena. This framework
466 encourages a research approach that is not centred on a single level, but instead integrates levels
467 through bottom-up (dark blue arrows) and top-down (light blue arrows) approaches.

468

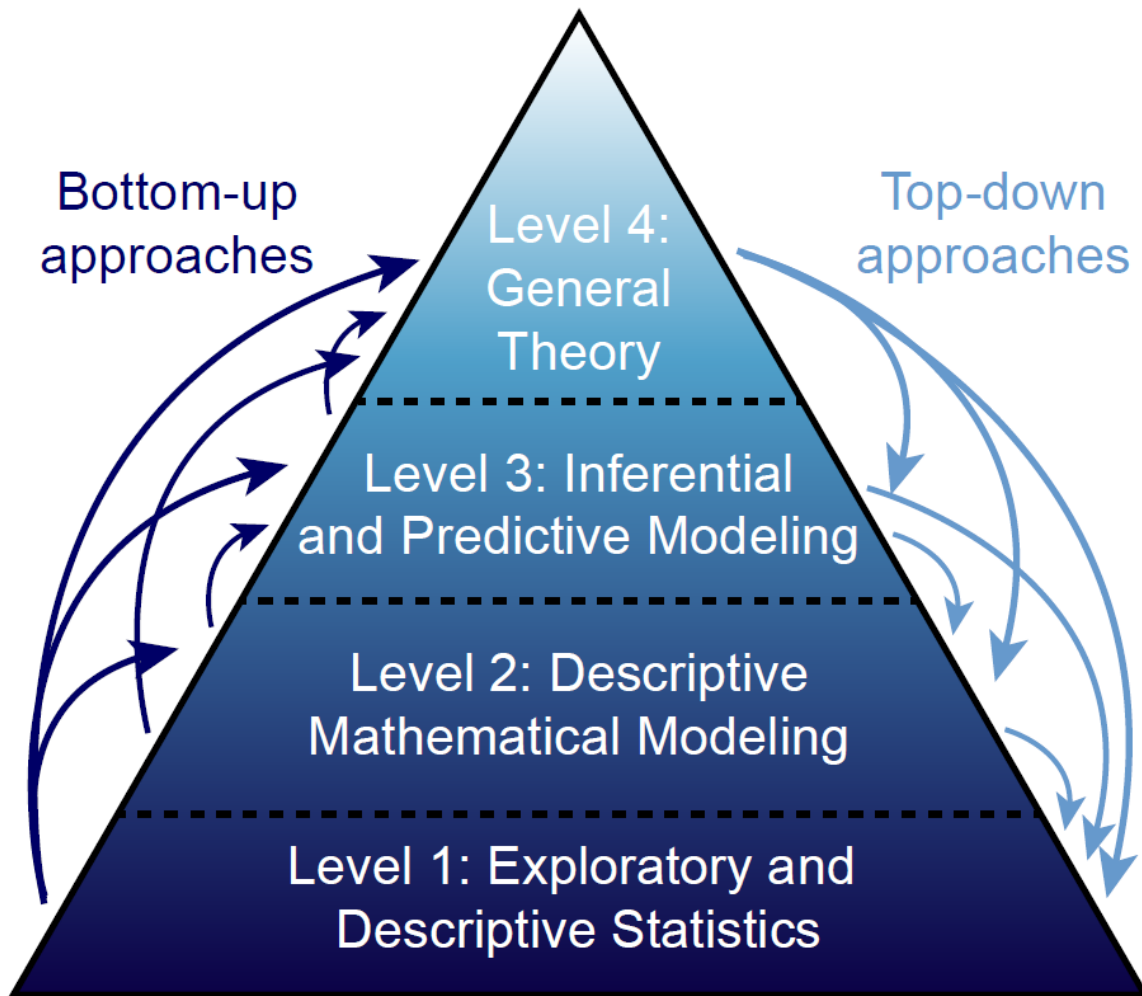
469 Figure 1

470



471

472



476 **Table 1. Key linguistic laws, including all that have been explicitly investigated beyond**
477 **human language.** For each law, we show the pair of variables related (Var. 1 and Var. 2),
478 mathematical models that have been proposed as a first approximation, and the sign of the
479 correlation between the variables when the law holds. A ‘type’ is an abstraction for a distinct unit
480 (e.g. a word or a behaviour), r is the frequency rank of a type (the most frequent has rank 1, the
481 2nd most frequent has rank 2, etc.), f is the frequency of a type, n is the number of types (that have
482 a certain frequency), l is the size of a type (e.g. its length or duration), μ is the number of meanings
483 or a proxy thereof (e.g. behavioural contexts of a type), t is the number of tokens, S_w is the size of
484 a whole construct (usually its number of parts) and S_p is the size of its parts. * is used to indicate
485 correlations that are inevitable given the definition of the variables involved. In the models, c
486 indicates a positive proportionality factor that, depending on the law, is not a free parameter
487 because it can be deduced before fitting (e.g. applying normalization). $\alpha, \beta, \gamma, \delta$ and η are positive
488 constants, the so-called exponents of the laws. a and b are additional parameters.

489 (For a more comprehensive bibliography on linguistic laws in biology, see:

490 <https://cqlab.upc.edu/biblio/laws/>)

491
492

Definitions of laws					Level of biological organisation at which the law has been explored, with examples of studies		
Law	Var. 1	Var. 2	Model	Correlation	Molecular	Organismal	Ecological
Zipf's rank-frequency law. This law is popularly known as Zipf's law but is only one of many laws in Zipf's popular book [4]. α is a positive parameter that is the so-called exponent of the law. $\alpha \approx 1$ for English words. By definition of rank, a negative correlation between f and r is expected for $\alpha > 0$.	r	f	$f = cr^\alpha$	$\leq 0^*$	Codons [108] DNA [9]	Vocal communication [42] Gestural communication [20]	Intra-population abundance distribution [28] Species abundance distributions [29] Disease spread [30]
Zipf-Mandelbrot law. This is a refinement of Zipf's rank-frequency law by Mandelbrot [109] that consists of introducing an additional parameter b . When $b=0$, one obtains Zipf's rank-frequency law.	r	r	$f = c(r+b)^\alpha$	$\leq 0^*$	RNA [10]	Vocal communication [21]	Intra-population abundance distribution [28] Species abundance distributions [31]
Zipf's number-frequency law. This law and Zipf's rank-frequency law are two sides of the same coin and the relationship between their respective exponents obeys approximately $\beta = 1/\alpha + 1$ [110]. Contrary to Zipf's rank-frequency law, a correlation between n and f is not expected <i>a priori</i> . β is a positive parameter that is the so-called exponent of this law. $\beta \approx 1$ for English words [96].	f	n	$n = cf^\beta$	≥ 0	Proteins [37]		
Zipf's law of abbreviation. Zipf found that more frequent words tend to be shorter [4] but he did not propose a function for the relationship between f and l . In a popular article, Sigurd et al [111] adopted an equation shown here that suggests that word frequency is determined by word	f	l	$f = cl^\beta$	≤ 0		Vocal communication [22] Gestural communication [23] Non-vocal behaviour [24]	

length while information theory suggests that it is rather the other way around [85]. See also [112].							
Zipf's law of meaning distribution. γ is the exponent of this law, that is $\gamma \approx 0.5$ for English words [113,114].	r		$\mu = Cr^\gamma$	≤ 0		Vocal communication [19]	
Zipf's meaning-frequency law. This law is a prediction by Zipf from the rank-frequency law and the law of meaning distribution [113]. He predicted $\delta \approx 0.5$ for English words, and later it was proven that $\delta = \gamma/\alpha$ [115].	f		$\mu = Cf^\delta$	≥ 0		Vocal communication [19] Gestural communication [25]	
Herdan's law (also known as Heaps' law). The law defines the growth of the number of distinct words as a function of the text length measured in tokens [67,116].	t	n	$n = Ct^\gamma$	$\geq 0^*$	Proteomes [13]		
Menzerath's law. The law is usually defined as a negative correlation between S_w and S_p . It bears the name of P. Menzerath [2,5], who was inspired by the relationship between the number of syllables of a word (S_w) and the duration of its syllables (S_p).	S_w	S_p	-	≤ 0		Vocal communication [26] Gestural communication [23]	
Menzerath-Altmann law. A generalization and mathematical formulation of Menzerath's law by G. Altmann [6]. In addition to the proportionality parameter c , it is defined by two additional parameters, a and b . a is an exponent that usually takes negative values. Although Solé [79] claimed that $a = -1$ and $b = 0$ are inevitable when S_p is defined as a mean size, the argument turned out to be flawed [98].	S_w	S_p	$S_p = cS_w^a e^{bS_w}$	≤ 0	Genes [14] Genomes [15] Proteins [16] RNA [17]		

493

494

495 **References**

- 496 1 Zipf, G.K. (1936) The psycho-biology of language: an introduction to dynamic philology.
497 George Routledge and Sons Ltd.
- 498 2 Menzerath, P. and De Oleza, J.M. (1928) *Spanische Lautdauer. Eine experimentelle*
499 *Untersuchung*, De Gruyter.
- 500 3 Köhler, R. *et al.* (2005) *Quantitative linguistics: An international handbook*, Walter de
501 Gruyter.
- 502 4 Zipf, G.K. (1949) *Human behavior and the principle of least effort*, Addison-Wesley Press.
- 503 5 Menzerath, P. (1954) *Die Architektonik des deutschen Wortschatzes*, Bonn.
- 504 6 Altmann, G. (1980) Prolegomena to Menzerath's law. *Glottometrika* 2, 1–10
- 505 7 Best, K.-H. and Rottmann, O. (2017) *Quantitative Linguistics: an Invitation, Studies in*
506 *Quantitative Linguistics* 25, RAM-Verlag.
- 507 8 Wilde, J. and Schwibbe, M.H. (1989) Organisationsformen von Erbinformation Im Hinblick
508 auf die Menzerathsche Regel. In *Das Menzerathsche Gesetz in informationsverarbeitenden*
509 *Systemen* (Altmann, G. *et al.*, eds), pp. 92–107, G. Olms
- 510 9 Mantegna, R.N. *et al.* (1994) Linguistic features of noncoding DNA sequences. *Phys. Rev.*
511 *Lett.* 73, 3169–3172
- 512 10 Schuster, P. *et al.* (1994) From sequences to shapes and back: A case study in RNA
513 secondary structures. *Proc. R. Soc. B Biol. Sci.* 255, 279–284
- 514 11 Huynen, M.A. and Van Nimwegen, E. (1998) The frequency distribution of gene family sizes
515 in complete genomes. *Mol. Biol. Evol.* 15, 583–589
- 516 12 Hoyle, D.C. *et al.* (2002) Making sense of microarray data distributions. *Bioinformatics* 18,
517 576–584
- 518 13 Nasir, A. *et al.* (2017) Phylogenetic tracings of proteome size support the gradual accretion
519 of protein structural domains and the early origin of viruses from primordial cells. *Front.*
520 *Microbiol.* 8, 1178
- 521 14 Nikolaou, C. (2014) Menzerath-Altman law in mammalian exons reflects the dynamics of
522 gene structure evolution. *Comput. Biol. Chem.* 53, 134–143
- 523 15 Ferrer-i-Cancho, R. and Forns, N. (2010) The self-organization of genomes. *Complexity* 15,
524 34–36
- 525 16 Shahzad, K. *et al.* (2015) The organization of domains in proteins obeys Menzerath-
526 Altmann's law of language. *BMC Syst. Biol.* 9, 44
- 527 17 Sun, F. and Caetano-Anollés, G. (2021) Menzerath–Altmann's law of syntax in RNA
528 accretion history. *Life* 11, 489
- 529 18 Calabrese, F. *et al.* (2019) Quantitation and comparison of phenotypic heterogeneity among
530 single cells of monoclonal microbial populations. *Front. Microbiol.* 10, 2814
- 531 19 Ferrer-i-Cancho, R. and McCowan, B. (2009) A law of word meaning in dolphin whistle
532 types. *Entropy* 11, 688–701
- 533 20 Genty, E. and Byrne, R.W. (2010) Why do gorillas make sequences of gestures? *Anim.*
534 *Cogn.* 13, 287–301
- 535 21 Hailman, J.P. (1994) Constrained permutation in “Chick-a-dee”- like calls of a black-lored tit
536 *Parus xanthogenys*. *Bioacoustics* 6, 33–50
- 537 22 Favaro, L. *et al.* (2020) Do penguins' vocal sequences conform to linguistic laws? *Biol. Lett.*

- 538 16,
- 539 23 Heesen, R. *et al.* (2019) Linguistic laws in chimpanzee gestural communication. *Proc. R. Soc. B Biol. Sci.* 286, 20182900
- 540
- 541 24 Ferrer-i-Cancho, R. and Lusseau, D. (2009) Efficient coding in dolphin surface behavioral
- 542 patterns. *Complexity* 14, 23–25
- 543 25 Hobaiter, C. and Byrne, R.W. (2014) The meanings of chimpanzee gestures. *Curr. Biol.* 24,
- 544 1596–1600
- 545 26 Gustison, M.L. *et al.* (2016) Gelada vocal sequences follow Menzerath’s linguistic law. *Proc. Natl. Acad. Sci.* 113, E2750–E2758
- 546
- 547 27 James, L.S. *et al.* (2021) Phylogeny and mechanisms of shared hierarchical patterns in
- 548 birdsong. *Curr. Biol.* 31, 2796–2808.
- 549 28 Kunakh, O.N. *et al.* (2018) Fitting competing models and evaluation of model parameters of
- 550 the abundance distribution of the land snail *Vallonia pulchella* (Pulmonata, Valloniidae).
- 551 *Regul. Mech. Biosyst.* 9, 198–202
- 552 29 Stedille, L.I.B. *et al.* (2020) Passive restoration in Araucaria Forest: useful ecological
- 553 indicators in monitoring successional advancement in exotic tree plantation landscapes.
- 554 *Restor. Ecol.* 28, 1213–1224
- 555 30 Kennedy, A.P. and Phillip Yam, S.C. (2020) On the authenticity of COVID-19 case figures.
- 556 *PLoS One* 15, e0243123
- 557 31 Guo, J. *et al.* (2020) Soil fungal assemblage complexity is dependent on soil fertility and
- 558 dominated by deterministic processes. *New Phytol.* 226, 232–243
- 559 32 Altmann, E.G. and Gerlach, M. (2016) Statistical laws in linguistics. In *Creativity and*
- 560 *Universality in Language. Lecture Notes in Morphogenesis.* (Degli, E. M. *et al.*, eds), pp. 7–
- 561 26, Springer, Cham
- 562 33 Bian, C. *et al.* (2016) Scaling laws and model of words organization in spoken and written
- 563 language. *EPL* 113, 18002
- 564 34 Balasubrahmanyam, V.K. and Naranan, S. (1996) Quantitative linguistics and complex
- 565 system studies. *J. Quant. Linguist.* 3, 177–228
- 566 35 Mehri, A. and Jamaati, M. (2017) Variation of Zipf’s exponent in one hundred live languages:
- 567 A study of the Holy Bible translations. *Phys. Lett. A* 381, 2470–2477
- 568 36 Furusawa, C. and Kaneko, K. (2003) Zipf’s law in gene expression. *Phys. Rev. Lett.* 90,
- 569 088102
- 570 37 Qian, J. *et al.* (2001) Protein family and fold occurrence in genomes: Power-law behaviour
- 571 and evolutionary model. *J. Mol. Biol.* 313, 673–681
- 572 38 Hailman, J.P. *et al.* (1985) The ‘chick-a-dee’ calls of *Parus atricapillus*: A recombinant
- 573 system of animal communication compared with written English. *Semiotica* 56, 191–224
- 574 39 Ficken, M.S. *et al.* (1994) The chick-a-dee call system of the Mexican chickadee. *Condor*
- 575 96, 70–82
- 576 40 Freeberg, T.M. and Lucas, J.R. (2012) Information theoretical approaches to chick-a-dee
- 577 calls of Carolina chickadees (*Poecile carolinensis*). *J. Comp. Psychol.* 126, 68–81
- 578 41 Markov, V.I. and Ostrovskaya, V.M. (1990) Organization of communication system in
- 579 *Tursiops truncatus montagu*. In *Sensory Abilities of Cetaceans* pp. 599–622, Springer US
- 580 42 McCowan, B. *et al.* (2002) Using information theory to assess the diversity, complexity, and
- 581 development of communicative repertoires. *J. Comp. Psychol.* 116, 166–172

- 582 43 Kershenbaum, A. *et al.* (2021) Shannon entropy as a robust estimator of Zipf's law in animal
583 vocal communication repertoires. *Methods Ecol. Evol.* 12, 553–564
- 584 44 Su, Q. (2018) A general pattern of the species abundance distribution. *PeerJ* 2018,
- 585 45 Bentz, C. and Ferrer-I-Cancho, R. (2016) Zipf's law of abbreviation as a language universal.
586 *Proc. Leiden Work. Capturing Phylogenetic Algorithms Linguist.* at
587 <<https://publikationen.uni-tuebingen.de/xmlui/handle/10900/68558>.>
- 588 46 Gahl, S. (2008) Time and thyme are not homophones: The effect of lemma frequency on
589 word durations in spontaneous speech. *Language (Baltim)*. 84, 474–496
- 590 47 Hernández-Fernández, A. *et al.* (2019) Linguistic laws in speech: The case of Catalan and
591 Spanish. *Entropy* 21, 1153
- 592 48 Torre, I.G. *et al.* (2019) On the physical origin of linguistic laws and lognormality in speech.
593 *R. Soc. Open Sci.* 6, 191023
- 594 49 Börstell, C. *et al.* (2016) Distribution and duration of signs and parts of speech in Swedish
595 Sign Language. *Sign Lang. Linguist.* 19, 143–196
- 596 50 Ficken, M.S. *et al.* (1978) A model of repetitive behaviour illustrated by chickadee calling.
597 *Anim. Behav.* 26, 630–633
- 598 51 Semple, S. *et al.* (2010) Efficiency of coding in macaque vocal communication. *Biol. Lett.* 6,
599 469–71
- 600 52 Valente, D. *et al.* (2021) Linguistic laws of brevity: conformity in *Indri indri*. *Anim. Cogn.* 24,
601 897–906
- 602 53 Ferrer-i-Cancho, R. and Hernández-Fernández, A. (2013) The failure of the law of brevity in
603 two new world primates. *Statistical caveats. Glottotheory* 4, 45–55
- 604 54 Bezerra, B.M. *et al.* (2011) Brevity is not always a virtue in primate communication. *Biol.*
605 *Lett.* 7, 23–5
- 606 55 White, E.P. *et al.* (2007) Relationships between body size and abundance in ecology.
607 *Trends Ecol. Evol.* 22, 323–330
- 608 56 Cramer, I. (2005) The parameters of the Altmann-Menzerath law. *J. Quant. Linguist.* 12, 41–
609 52
- 610 57 Milička, J. (2014) Menzerath's Law: The whole is greater than the sum of its parts. *J. Quant.*
611 *Linguist.* 21, 85–99
- 612 58 Andres, J. *et al.* (2019) Towards a fractal analysis of the sign language. *J. Quant. Linguist.*
613 DOI: 10.1080/09296174.2019.1656149
- 614 59 Hou, R. *et al.* (2017) A study on correlation between Chinese sentence and constituting
615 clauses based on the Menzerath-Altmann law. *J. Quant. Linguist.* 24, 350–366
- 616 60 Torre, I.G. *et al.* (In Press) Can Menzerath's law be a criterion of complexity in
617 communication? *PLoS One*
- 618 61 Grzybek, P. and Köhler, R. (2007) Do we have problems with Arens' law? A new look at the
619 sentence-word relation. In *Exact Methods in the Study of Language and Text: Dedicated to*
620 *Gabriel Altmann on the Occasion of his 75th Birthday. Quantitative Linguistics*, 62. Mouton
621 de Gruyter, 205–217
- 622 62 Li, W. (2012) Menzerath's law at the gene-exon level in the human genome. *Complexity* 17,
623 49–53
- 624 63 Gustison, M.L. and Bergman, T.J. (2017) Divergent acoustic properties of gelada and
625 baboon vocalizations and their implications for the evolution of human speech. *J. Lang.*

- 626 *Evol.* 2, 20–36
- 627 64 Fedurek, P. *et al.* (2017) Trade-offs in the production of animal vocal sequences: Insights
628 from the structure of wild chimpanzee pant hoots. *Front. Zool.* 14, 50
- 629 65 Huang, M. *et al.* (2020) Male gibbon loud morning calls conform to Zipf’s law of brevity and
630 Menzerath’s law: insights into the origin of human language. *Anim. Behav.* 160, 145–155
- 631 66 Clink, D.J. and Lau, A.R. (2020) Adherence to Menzerath’s law is the exception (not the
632 rule) in three duetting primate species. *R. Soc. Open Sci.* 7, 201557
- 633 67 Herdan, G. (1960) *Type-token mathematics*, Mouton.
- 634 68 Demartsev, V. *et al.* (2019) The “Law of Brevity” in animal communication: Sex-specific
635 signaling optimization is determined by call amplitude rather than duration. *Evol. Lett.* 3,
636 623–634
- 637 69 Evans, N. and Levinson, S.C. (2009) The myth of language universals: language diversity
638 and its importance for cognitive science. *Behav. Brain Sci.* 32, 429–92
- 639 70 Clink, D.J. *et al.* (2020) Brevity is not a universal in animal communication: evidence for
640 compression depends on the unit of analysis in small ape vocalizations. *R. Soc. Open Sci.*
641 7, 200151
- 642 71 Luo, B. *et al.* (2013) Brevity is prevalent in bat short-range communication. *J. Comp.*
643 *Physiol. A. Neuroethol. Sens. Neural. Behav. Physiol.* 199, 325–33
- 644 72 Watson, S.K. *et al.* (2020) An exploration of Menzerath’s law in wild mountain gorillas. *Open*
645 *Sci. Fram.* 16, 20200380
- 646 73 Altmann, G. (1993) Science and linguistics. In *Contributions to Quantitative Linguistics*
647 (Köhler, R. and Rieger, B., eds), pp. 3–10, Kluwer
- 648 74 Bunge, M. (2001) *La science, sa méthode et sa philosophie*, Vigdor.
- 649 75 Miller, G.A. (1957) Some effects of intermittent silence. *Am. J. Psychol.* 70, 311–314
- 650 76 Howes, D. (1968) Zipf’s law and Miller’s random-monkey model. *Am. J. Psychol.* 81, 269–
651 272
- 652 77 Simon, H.A. (1955) On a class of skew distribution functions. *Biometrika* 42, 425
- 653 78 Li, W. (1992) Random texts exhibit Zipf’s-law-like word frequency distribution. *IEEE Trans.*
654 *Inf. Theory* 38, 1842–1845
- 655 79 Solé, R. V (2010) Genome size, self-organization and DNA’s dark matter. *Complexity* 16,
656 20–23
- 657 80 Ferrer-i-Cancho, R. and Elvevåg, B. (2010) Random texts do not exhibit the real Zipf’s law-
658 like rank distribution. *PLoS One* 5, 29411
- 659 81 Ferrer-i-Cancho, R. *et al.* (2012) The challenges of statistical patterns of language: The case
660 of Menzerath’s law in genomes. *Complexity* 18, 11–17
- 661 82 Carrera-Casado, D. and Ferrer-i-Cancho, R. (In Press) The advent and fall of a vocabulary
662 learning bias from communicative efficiency. *Biosemiotics* at
663 <<http://arxiv.org/abs/2105.11519>>
- 664 83 Plotkin, J.B. and Nowak, M.A. (2000) Language evolution and information theory. *J. Theor.*
665 *Biol.* 205, 147–159
- 666 84 Ferrer-i-Cancho, R. *et al.* (2013) Compression as a universal principle of animal behavior.
667 *Cogn. Sci.* 37, 1565–1578
- 668 85 Ferrer-i-Cancho, R. *et al.* (2020) Optimal coding and the origins of Zipfian laws. *J. Quant.*

- 669 *Linguist.* DOI: 10.1080/09296174.2020.1778387
- 670 86 Ferrer-i-Cancho, R. (2016) Compression and the origins of Zipf's law for word frequencies.
671 *Complexity* 21, 409–411
- 672 87 Stumpf, M.P.H. and Porter, M.A. (2012) Critical truths about power laws. *Science*. 335, 665–
673 666
- 674 88 Ferrer-i-Cancho, R. (2018) Optimization models of natural communication. *J. Quant.*
675 *Linguist.* 25, 207–237
- 676 89 Caetano-Anollés, G. (2021) The compressed vocabulary of microbial life. *Front. Microbiol.*
677 12, 655990
- 678 90 Robert Burger, J. *et al.* (2021) Universal rules of life: metabolic rates, biological times and
679 the equal fitness paradigm. *Ecol. Lett.* 24, 1262–1281
- 680 91 Brown, J.H. *et al.* (2004) Toward a metabolic theory of ecology. *Ecology* 85, 1771–1789
- 681 92 Brown, J.H. *et al.* (2018) Equal fitness paradigm explained by a trade-off between
682 generation time and energy production rate. *Nat. Ecol. Evol.* 2, 262–268
- 683 93 Gerlach, M. and Altmann, E.G. (2019) Testing statistical laws in complex systems. *Phys.*
684 *Rev. Lett.* 122, 168301
- 685 94 Li, W. *et al.* (2010) Fitting ranked linguistic data with two-parameter functions. *Entropy* 12,
686 1743–1764
- 687 95 Font-Clos, F. and Corral, Á. (2015) Log-log convexity of type-token growth in Zipf's systems.
688 *Phys. Rev. Lett.* 114, 238701
- 689 96 Moreno-Sánchez, I. *et al.* (2016) Large-scale analysis of Zipf's law in English texts. *PLoS*
690 *One* 11, e0147073
- 691 97 Corral, Á. *et al.* (2020) Distinct flavors of Zipf's law and its maximum likelihood fitting: Rank-
692 size and size-distribution representations. *Phys. Rev. E* 102, 052113
- 693 98 Ferrer-i-Cancho, R. *et al.* (2014) When is Menzerath-Altmann law mathematically trivial? A
694 new approach. *Stat. Appl. Genet. Mol. Biol.* 13, 633–644
- 695 99 Semple, S. *et al.* (2013) The Law of Brevity in macaque vocal communication is not an
696 artefact of analysing mean call durations. *J. Quant. Linguist.* 20, 209–217
- 697 100 Deluca, A. and Corral, Á. (2013) Fitting and goodness-of-fit test of non-truncated and
698 truncated power-law distributions. *Acta Geophys.* 61, 1351–1394
- 699 101 Miller, G.A. and Chomsky, N. (1963) Finitary models for language users. In *Handbook of*
700 *Mathematical Psycholgy* (2nd edn) (Luce, D., ed), pp. 419–492, John Wiley & Sons
- 701 102 Bonhoeffer, S. *et al.* (1996) No signs of hidden language in noncoding DNA. *Phys. Rev.*
702 *Lett.* 76, 1977
- 703 103 Miton, H. and Morin, O. (2019) When iconicity stands in the way of abbreviation: No Zipfian
704 effect for figurative signals. *PLoS One* 14, e0220793
- 705 104 Ferrer-i-Cancho, R. and Gavaldà, R. (2009) The frequency spectrum of finite samples from
706 intermittent silence process. *J. Am. Soc. Inf. Sci. Technol.* 64, 837–843
- 707 105 Ferrer-i-Cancho, R. and McCowan, B. (2012) The span of correlations in dolphin whistle
708 sequences. *J. Stat. Mech. Theory Exp.* 2012, P06002
- 709 106 Piantadosi, S.T. (2014) Zipf's word frequency law in natural language: A critical review and
710 future directions. *Psychon. Bull. Rev.* 21, 1112–1130
- 711 107 Suzuki, R. *et al.* (2005) The use of Zipf's law in animal communication analysis. *Anim.*

- 712 *Behav.* 69, F9–F17
- 713 108 Balasubrahmanyam, V.K. and Narayan, S. (2000) Information theory and algorithmic
714 complexity: Applications to language discourses and DNA sequences as complex systems
715 Part II: Complexity of DNA sequences, analogy with linguistic discourses. *J. Quant. Linguist.*
716 7, 153–183
- 717 109 Mandelbrot, B. (1954) Structure formelle des textes et communication. *WORD* 10, 1–27
- 718 110 Pietronero, L. *et al.* (2001) Explaining the uneven distribution of numbers in nature: The laws
719 of Benford and Zipf. *Phys. A Stat. Mech. its Appl.* 293, 297–304
- 720 111 Sigurd, B. *et al.* (2004) Word length, sentence length and frequency - Zipf revisited. *Stud.*
721 *Linguist.* 58, 37–52
- 722 112 Corral, Á. and Serra, I. (2020) The brevity law as a scaling law, and a possible origin of
723 Zipf's law for word frequencies. *Entropy* 22, 1–14
- 724 113 Zipf, G.K. (1945) The meaning-frequency relationship of words. *J. Gen. Psychol.* 33, 251–
725 256
- 726 114 Casas, B. *et al.* (2019) Polysemy and brevity versus frequency in language. *Comput.*
727 *Speech Lang.* 58, 19–50
- 728 115 Ferrer-i-Cancho, R. and Vitevitch, M.S. (2018) The origins of Zipf's meaning-frequency law.
729 *J. Assoc. Inf. Sci. Technol.* 69, 1369–1379
- 730 116 Heaps, H.S. (1978) *Information Retrieval: Computational and Theoretical Aspects*,
731 Academic Press.
- 732

733 **Highlights**

734

735 • Linguistic laws refer to statistical patterns shared across human languages. Investigation of these
736 patterns has been extended to a range of biological systems, from molecules to organisms to
737 ecosystems, with the number of studies increasing in recent years.

738 • Linguistic laws and established concepts in different fields of biology may describe similar - or
739 even identical - patterns, providing an opportunity for unification of natural and language sciences.

740 • We propose an overarching framework which shifts the emphasis from exploring and describing
741 linguistic laws, to identifying underlying mechanisms, generating predictions, and ultimately
742 developing general theory about the organization of biological systems.

743 • This potential to develop new theory for understanding the natural world will only be realised
744 through cross-fertilisation of ideas between researchers working in diverse disciplines and
745 focussed on different levels of biological organisation.

746

747 **Outstanding Questions**

748 • **Which linguistic laws hold in which biological systems?**

749 It is vital to explore whether an absence of evidence for linguistics laws in different systems reflects
750 reality or just a lack of exploration; critically, many linguistic laws have never been explicitly tested
751 beyond language, and investigating these could provide novel insights.

752
753 • **What underpins conformity to linguistic laws in different biological systems?**

754 Broad scale comparative approaches will allow testing of hypotheses about the functionality and
755 potential ecological and/or evolutionary drivers of patterns consistent with linguistic laws, and
756 investigation of whether patterns reflect the sharing of a common ancestral trait (constraint) or
757 convergent evolution (selection).

758
759 • **What are the most appropriate currencies to use when exploring linguistic laws in
760 biological systems?**

761 Comparative approaches necessitate the development of measures for testing laws that are
762 applicable across taxa in one system, but also across different systems and levels of biological
763 organization (for example, direct or indirect measures of energetic cost).

764
765 • **How does the manifestation of linguistic laws shift across time in biological systems?**

766 Temporal variation in conformity to linguistic laws is poorly understood; it is important now to
767 explore how these patterns are shaped by both internal processes (e.g., physiological states,
768 ontogenesis) and external conditions (e.g., environmental contexts).

769
770 • **Should we rename linguistic laws, and if so how?**

771 As linguistic laws hold in diverse biological systems, arguably they should be renamed to reflect
772 this broader applicability; any change in terminology should consider the commonality between
773 laws and key concepts in different fields.

774
775 • **What are key areas and aims for future studies of linguistic laws in biology?**

776 Investigations of linguistic laws at levels of biological organisation (e.g. cells, tissues, organs) and
777 in disciplinary fields (e.g. biosemiotics, developmental biology) where they have rarely – or never –
778 been explored (Figure 1) could open up new opportunities for cross-disciplinary integration, and
779 further contribute to development of general theory (Figure 2).

780