

1 DNA as a Language: A Linguistic Comparison with

2 Human Spoken and Written Systems

3 Jack Cohen* and Barak Akabayov

4 Department of Chemistry and Data Science Research Center, Ben-Gurion University of the Negev,
5 Beer-Sheva 8410501, Israel

6

7 **Received:** Feb 08, 2026 **Revised:** Jun 04, 2026 **Just Accepted Online:** Jun 11, 2026
8 **Published:** Xxx

9 This article has been accepted for publication and undergone full peer review but has not been
10 through the copyediting, typesetting, pagination and proofreading process, which may lead to
11 differences between this version and the Version of Record.

12 Please cite this article as:

13 J. Cohen, B. Akabayov (2026) DNA as a Language: A Linguistic Comparison with Human
14 Spoken and Written Systems. **Substantia**. *Just Accepted*. DOI: 10.36253/Substantia-3980

15

16 Abstract

17 The traditional metaphor of DNA as a language is examined through systematic analysis,
18 using formal linguistic categories rather than relying on heuristic analogies. By exploring
19 aspects of phonology, morphology, syntax, semantics, pragmatics, and diachrony, the study
20 investigates the structural parallels between genetic sequences and human spoken and written
21 languages, while also identifying the fundamental limitations of this comparison. The
22 analysis indicates that DNA displays distinct symbolic units, a combinatorial organization,
23 strict ordering constraints, context dependence, and historical evolution. These characteristics
24 resemble written language more closely than they do speech. However, DNA fundamentally
25 differs from human language in several ways: it lacks intentionality, reference, generativity,

26 and representational meaning. Genetic “meaning” is defined by its operational and causal
27 roles, focusing on biochemical functions rather than semantic interpretations. This study
28 emphasizes that DNA should not be viewed as a literal language but rather as a system of
29 symbols grounded in chemistry. The comparison with linguistics is most useful for
30 establishing distinctions rather than for expanding metaphorical interpretations.

31 **Keywords:** DNA language, morphemes, codons, gene expression, linguistics

32 The metaphor of DNA as a “language” has circulated in molecular biology since the mid-
33 twentieth century, often as a heuristic for explaining how genetic sequences encode and
34 transmit information. When examined through the formal methodology of linguistics rather
35 than metaphor alone, this comparison becomes both more precise and more constrained.
36 Linguistics offers a multi-level analytical framework—encompassing phonology,
37 morphology, syntax, semantics, pragmatics, and diachrony—that allows for systematic
38 comparison between symbol systems. Applying these categories to DNA clarifies where the
39 analogy with human spoken and written languages is structurally informative and where it
40 breaks down in principle.

41 **Units of Expression: Phonology and Orthography**

42 Human spoken languages are organized around phonology: a finite inventory of discrete,
43 contrastive sound units (phonemes) that are meaningless in isolation but meaningful in
44 combination [1]. Written language introduces an orthographic layer, in which visual symbols
45 represent phonological units or morphemes depending on the writing system [2].

46 DNA lacks phonology in the acoustic sense, but it does possess a discrete symbolic inventory
47 analogous to an orthography. The four nucleotide bases—adenine, thymine, cytosine, and
48 guanine—form a finite set of contrastive units whose sequence determines downstream

49 effects in protein synthesis [3]. Unlike phonemes, however, nucleotides are not abstract
50 cognitive categories; they are chemically instantiated entities whose physical properties are
51 inseparable from their informational role [4].

52 In this respect, DNA more closely resembles written language than speech, yet the analogy
53 remains partial. Orthographic symbols conventionally stand for sounds or meanings external
54 to the script, whereas DNA bases directly participate in the biochemical processes that “read”
55 them. The symbol and its material realization cannot be cleanly separated, a point
56 emphasized in discussions of the symbol–matter relationship in biology [5]. For instance, in
57 the human β -globin gene, a single nucleotide substitution (A→T) changes the DNA codon
58 from GAG to GTG. This alteration results in the incorporation of a different amino acid—
59 valine instead of glutamic acid—into the hemoglobin protein, leading to the development of
60 sickle-cell anemia [6]. This example illustrates that, much as changing a single letter in a
61 word can significantly alter the “meaning” of a word, altering a single nucleotide can
62 significantly change the “meaning” in the language of genetics.

63 **Morphology: From Morphemes to Codons and Genes**

64 Morphology in linguistics concerns the internal structure of words and the smallest
65 meaning-bearing units, or morphemes. These units combine according to systematic rules and
66 often productively generate new lexical items [7].

67 DNA exhibits a comparable intermediate level of organization. Codons—triplets of
68 nucleotides—function as minimal functional units, each specifying an amino acid or a
69 termination signal in the genetic code [8]. In this sense, codons resemble morphemes more
70 closely than letters, as they map systematically to biochemical outputs. Genes, in turn,

71 resemble complex lexical items or fixed expressions, with internal organization that includes
72 coding regions, non-coding regions, and regulatory elements [4].

73 The crucial divergence lies in productivity and agency. Human speakers can generate novel
74 morphological forms intentionally and interpret them on the fly. Genetic systems cannot
75 invent new codon–amino acid correspondences or morphological rules; innovation occurs
76 only through mutation, recombination, and selection, processes that are blind rather than
77 generative in the linguistic sense [9]. For example, the codon "ATG" in DNA corresponds to
78 the amino acid methionine and acts as a universal "start" signal for protein synthesis. This
79 codon can be considered a morpheme-like unit of meaning in DNA because its presence at
80 the beginning of a gene indicates the initiation of translation. A complete gene, such as the
81 one encoding DNA primase, can be viewed as a larger lexical unit. It consists of coding
82 sequences (exons) that code for the primase protein, as well as regulatory sequences (such as
83 promoters and enhancers) that control when and where primase is produced. This is
84 analogous to a complex word formed by a root and affixes that modify its usage.

85 **Syntax: Linear Order and Structural Constraints**

86 Syntax governs how words combine into larger structures, imposing both linear and
87 hierarchical constraints. Human syntax is recursive and generative, enabling an unbounded
88 range of expressions from finite resources [10].

89 DNA also exhibits syntactic properties in a limited but non-trivial sense. Linear order is
90 critical: frame shifts, misplaced stop codons, or altered regulatory sequences can radically
91 change or eliminate function [3]. There are positional constraints analogous to syntactic
92 rules—promoters must precede genes, splice sites must occur at specific boundaries, and
93 regulatory motifs interact across defined spatial relationships [4].

94 However, DNA syntax is procedural rather than representational. It lacks recursion,
95 hierarchical embedding, and ambiguity resolution. Genetic sequences do not express
96 propositions; they trigger biochemical processes. The “grammar” of DNA specifies what
97 reactions occur, not what states of affairs are described. For example, when a single
98 nucleotide is deleted from a coding sequence, it shifts the reading frame of the triplet codons,
99 resulting in a nonsensical protein sequence or a premature stop. This type of frameshift
100 mutation is observed in the human ABO blood group gene. For example, a one-base deletion
101 (ΔG at position 261) in exon 6 alters the reading frame and produces a truncated,
102 nonfunctional enzyme, which is responsible for the O blood type [11]. Similarly, specific
103 sequence motifs—akin to grammar—are essential for proper DNA function. The bacterial
104 primase enzyme, for instance, does not initiate RNA primer synthesis at random DNA
105 sequences. *E. coli* primase (DnaG) ignores about 97% of possible trinucleotide sites and
106 initiates primers only at specific recognition sequences roughly every 1.5 to 2 kilobases apart
107 [12]. This process is analogous to the necessity of placing certain words or punctuation
108 correctly in a sentence for it to make sense. A misplaced stop codon, for instance, is like
109 inserting a period in the middle of a sentence; it prematurely ends the message. In DNA, an
110 incorrectly positioned stop codon aborts the protein synthesis, adhering to the "syntax" rules
111 of the genetic code.

112 **Semantics: Meaning Without Reference**

113 In linguistics, semantics concerns meaning and reference: how linguistic expressions relate to
114 concepts, entities, and states of affairs. Human language is intentional, referential, and
115 capable of truth or falsity [13].

116 DNA semantics operates under a fundamentally different logic. The “meaning” of a genetic
117 sequence is exhausted by its functional consequences within a cellular context—protein

118 structure, expression timing, or regulatory interaction [14]. A gene does not refer to a protein
119 in the way a word refers to an object or concept; it causally participates in its production. This
120 distinction has been emphasized in philosophical critiques of biological information that
121 caution against importing semantic notions wholesale from linguistics into biology [15].

122 Accordingly, DNA meaning cannot be false, metaphorical, or hypothetical. A sequence can
123 only function properly, malfunction, or remain unexpressed. This marks a categorical
124 difference between representational language and operational biological coding systems. For
125 example, the DNA sequence of a gene is considered "meaningful" only to the extent that it
126 produces a functional product. For example, the gene that encodes DNA primase doesn't
127 merely describe the primase; it actually serves as a template for assembling the primase
128 enzyme. If a mutation introduces a premature stop codon in that gene, the result is a
129 truncated, nonfunctional enzyme, which can lead to certain genetic diseases. However, we
130 wouldn't label the message of the mutated gene as "false"; we would simply describe it as
131 nonfunctional. In human language, a sentence can be false or nonsensical, but in DNA, a
132 "nonsensical" sequence merely fails to produce a viable outcome. For instance, when a codon
133 changes from a sense codon to a nonsense (stop) codon, it doesn't carry metaphorical
134 meaning; instead, it halts protein synthesis entirely, eliminating the gene's effect rather than
135 misrepresenting it.

136 **Pragmatics: Context Without Intention**

137 Pragmatics studies how context shapes meaning in language use, including speaker intention,
138 implicature, and social norms. While DNA lacks speakers and intentions, context nonetheless
139 plays a decisive role in genetic expression. The same sequence can yield different outcomes
140 depending on cellular environment, epigenetic modification, developmental stage, or external
141 conditions [4].

142 This context sensitivity is sometimes described in pragmatic terms within biosemiotics [16],
143 but the analogy remains limited. DNA interpretation involves no inference, negotiation, or
144 communicative goals. Context alters outcomes mechanically rather than conversationally. For
145 instance, the lac operon in *E. coli* contains genes responsible for lactose metabolism and is
146 expressed only when lactose is present and glucose is absent in the cell's environment. In any
147 other context, these genes—though present in the DNA—remain unexpressed or "silent." The
148 DNA sequence itself does not change; rather, it is the cellular context (the availability of
149 lactose and the absence of glucose) that enables the sequence to be read and acted upon.
150 Similarly, the gene for DNA primase is active only during the S-phase of the cell cycle, when
151 DNA replication occurs. The presence of specific replication signals and proteins, such as
152 helicase unwinding the DNA, provides the context in which the primase gene is transcribed
153 and its protein product utilized. Like a sentence that only makes sense in a particular
154 situational context, a genetic sequence's effect relies on biochemical context—albeit without
155 any conscious intention behind it.

156 **Diachrony: Evolution and Language Change**

157 Historical linguistics examines how languages change over time through sound shifts,
158 semantic drift, grammaticalization, and contact [17]. Genetic systems likewise change
159 through mutation, drift, recombination, and natural selection [8].

160 The parallel is strongest at this level. Both systems exhibit descent with modification,
161 conservation of core structures, and divergence into families. Yet the driving forces differ
162 fundamentally. Language change is mediated by cognition, culture, and social interaction
163 [18], whereas genetic change is governed by differential survival and reproduction. For
164 example, the DNA primase enzyme in bacteria and the primase found in humans have
165 diverged significantly in both sequence and subunit structure over billions of years. However,

166 they still share the core function of synthesizing RNA primers. This situation is similar to the
167 evolution of languages from Latin; for example, the Latin word "mater" evolved into "madre"
168 in Spanish and "mère" in French. Although these words sound different, they retain the same
169 meaning: "mother." In biology, there are distinct families of primases that exhibit no obvious
170 sequence similarity, suggesting they have taken independent evolutionary paths.
171 Nevertheless, even these divergent primases often share essential catalytic features due to the
172 functional constraints of their roles [19]. In both linguistic and biological contexts, the
173 process of inheritance with variation produces families—language families and gene
174 families—that preserve a core structure or meaning while diversifying over time.

175 **Conclusion**

176 Applying linguistic methodology to DNA reveals a system that is language-like in structure
177 but not linguistic in kind. DNA employs discrete symbols, combinatorial organization,
178 constrained order, context sensitivity, and historical evolution. At the same time, it lacks
179 intentionality, reference, productivity, and representation—the defining properties of human
180 spoken and written languages.

181 The value of the comparison lies not in treating DNA as a literal language, but in using
182 linguistic categories to specify precisely where the analogy holds and where it fails. DNA is
183 best understood as a chemically instantiated, operational symbol system—one that converges
184 with language at the level of structure while diverging decisively at the levels of meaning,
185 use, and cognition.

186 **References**

- 187 1. Hockett, C., *The Origin of Speech*. Scientific American, 1960. **203**: p. 88-96.
- 188 2. Saussure, F.d., *Course in General Linguistics* ed. T. W. Baskin. 1916/2011: Columbia
189 University Press.

- 190 3. Crick, F.H.C., *On protein synthesis*. Symposia of the Society for Experimental
191 Biology, 1958. **12**: p. 138–163.
- 192 4. Alberts, B., et al., *Molecular Biology of the Cell* 7th ed. ed. 2022: Garland Science.
- 193 5. Pattee, H.H., *Cell psychology: An evolutionary approach to the symbol–matter*
194 *problem*. Cognition and Brain Theory, 1982. **5**: p. 325–341.
- 195 6. Inusa, B.P.D., et al., *Sickle Cell Disease-Genetics, Pathophysiology, Clinical*
196 *Presentation and Treatment*. Int J Neonatal Screen, 2019. **5**(2): p. 20.
- 197 7. Jackendoff, R., *Foundations of Language (Brain, Meaning, Grammar, Evolution)*. 1st
198 ed. 2003: Oxford University Press. 477.
- 199 8. Crick, F.H., *The origin of the genetic code*. J Mol Biol, 1968. **38**(3): p. 367-79.
- 200 9. Maynard Smith, J., *The concept of information in biology*. Philosophy of Science,
201 2000. **67**(2): p. 177–194.
- 202 10. Chomsky, N., *Aspects of the Theory of Syntax*. . 1965: MIT Press. 251.
- 203 11. Kitano, T., A. Blancher, and N. Saitou, *The functional A allele was resurrected via*
204 *recombination in the human ABO blood group gene*. Mol Biol Evol, 2012. **29**(7): p.
205 1791-6.
- 206 12. Soffer, A., et al., *Inferring primase-DNA specific recognition using a data driven*
207 *approach*. Nucleic Acids Res, 2021. **49**(20): p. 11447-11458.
- 208 13. Saussure, F.d.W.B., Trans.). Columbia University Press., *Course in General*
209 *Linguistics (W. Baskin, Trans.)*. 1916/2011: Columbia University Press.
- 210 14. Barbieri, M., *The Organic Codes: An Introduction to Semantic Biology*. . 2003:
211 Cambridge University Press.
- 212 15. Godfrey-Smith, P., *Information, arbitrariness, and selection: Comments on Maynard*
213 *Smith*. . Philosophy of Science, 2000. **67**(2): p. 202–207.
- 214 16. Hoffmeyer, J., *Biosemiotics: An Examination into the Signs of Life and the Life of*
215 *Signs*. 2008: University of Scranton Press.
- 216 17. Labov, W., *Principles of Linguistic Change, Volume 1: Internal Factors*. . 1994:
217 Blackwell.
- 218 18. Croft, W., *Explaining Language Change: An Evolutionary Approach*. 2000:
219 Longman.
- 220 19. Kuchta, R.D. and G. Stengel, *Mechanism and evolution of DNA primases*. Biochim
221 Biophys Acta, 2010. **1804**(5): p. 1180-9.