

1 **Structural Enzymology, Phylogenetics, Differentiation, and Symbolic Reflexivity at the Dawn of**
2 **Biology**

3
4 Charles W. Carter, Jr^{1*}, Guo Qing Tang¹, Sourav Kumar Patra¹, Laurie Betts¹, Henry Dieckhaus^{1,2}, Brian
5 Kuhlman^{1,3}, Jordan Douglas^{4,5}, Peter R. Wills⁴, Remco Bouckaert⁵, Milena Popovic⁶, and Mark A. Ditzler⁶

6
7 ¹Department of Biochemistry and Biophysics, University of North Carolina at Chapel Hill, Chapel Hill, NC
8 27599-7260; ²Lineberger Comprehensive Cancer Center, School of Medicine, University of North Carolina
9 at Chapel Hill, Chapel Hill, NC ³Division of Chemical Biology and Medicinal Chemistry, Eshelman School
10 of Pharmacy, University of North Carolina at Chapel Hill, Chapel Hill, NC ⁴Department of Physics,
11 Auckland University, Auckland, NZ; ⁵Department of Computer Science, Auckland University, Auckland,
12 NZ; ⁶Blue Marble Space Institute of Science, Ames, CA; ⁷Center for the Emergence of Life, NASA Ames
13 Research Center, Moffett Field, CA

14
15 **Author for Correspondence*> Charles W. Carter, Jr. Department of Biochemistry and Biophysics,
16 University of North Carolina at Chapel Hill. Telephone: 919 259 2558. FAX: (919) 843-3328. Email:
17 carter@med.unc.edu

18 ORCID IDs:

19 C.W.C: 0000-0002-2653-4452

20 S.K.P: 0000-0002-9314-5641

21 L.B: 0000-0001-7594-3028

22 H.D: 0000-0003-1390-2444

23 B.K: 0000-0003-4907-9699

24 J.D: 0000-0003-0371-9961

25 P.R.W: 0000-0002-2670-7624

26 M.P: 0009-0007-9914-5558

27 M:D: 0000-0003-3108-4596

28 **Keywords:** Origin of genetic coding; Aminoacyl-tRNA synthetase•RNA cognate pairs; Bidirectional
29 genetic coding; ancestral genes; substrate specificity; phylogenetics

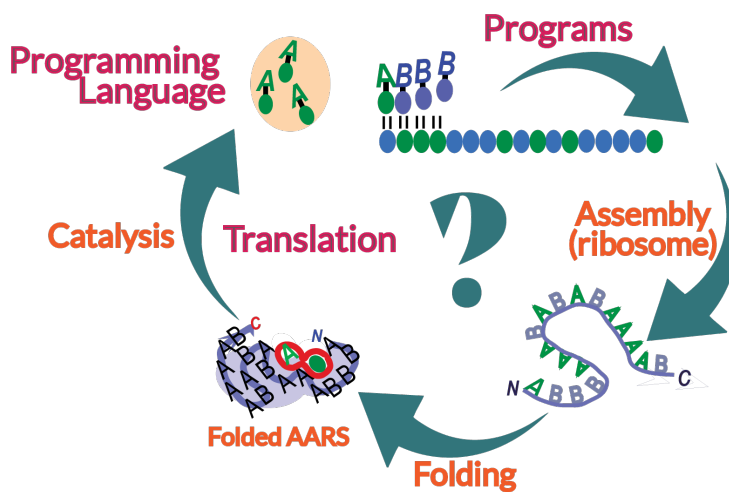
30

31 **Graphical Abstract**
32

Biology's Operating System

How did these happen?

Software | Hardware



33
34
35

36

37 **Abstract**

38 The reflexive translation of symbols in one chemical language to another defined genetics. Yet, the co-
39 linearity of codons and amino acids is so commonplace an idea that few even ask how it arose. Readout is
40 done by two distinct sets of proteins, called aminoacyl-tRNA synthetases (AARS). AARS must enforce the
41 rules first used to assemble themselves. The roots of translation lie in experimentally testing the structural
42 codes that the earliest AARS•tRNA cognate pairs used to recognize both amino acid and RNA substrates.
43 We review here new results on five different facets of that problem. (i) The surfaces of structures coded by
44 opposite strands of the same gene have opposite polarities. The corresponding proteins then fold up “inside
45 out” relative to one another. The inversion symmetry of base pairing thus projects into the proteome. That
46 leads in turn to contrasting amino acid and RNA substrate binding modes. (ii) *E. coli* reproduces *in vivo* the
47 nested hierarchy of active excerpts we had designed as models—protozymes and urzymes—for ancestral
48 AARS. (iii) A third novel deletion produced *in vivo* and a new Class II urzyme suggest how to design
49 bidirectional urzyme genes. (iv) Codon middle-base pairing provides a basis to constrain Class I and II
50 AARS family trees. (v) AARS urzymes acylate Class-specific subsets of an RNA library, showing RNA
51 substrate specificity for the first time. Four new phylogenetic routines augment these results to compose a
52 viable platform for experimental study of the origins of genetic coding.

53

54 Significance Statement

55 The origin of genetic coding poses questions distinct from those faced studying the evolution of
56 enzymes since the first cells. Modern enzymes that translate the code range in size from ~330 to ~970
57 amino acids. Ancestral forms cannot have been nearly as complex. Moreover, such primitive enzymes
58 likely could enforce only a much-reduced coding alphabet. Structural and molecular biology data
59 point to a broad sketch of events leading to the code. That research platform will enable us to see how
60 Nature came to store information about the physical chemistry of amino acids in the coding table.
61 That, in turn, allowed searching of a very broad amino acid sequence space. Selection could then
62 learn how to assemble amino acids into functional, reflexive catalysts. Those catalysts had rates and
63 fidelities consistent with bootstrapping the modern coding alphabet. New phylogenetic algorithms
64 need to be developed to fully test that putative sketch experimentally.

65

66 Introduction

67 Genetic coding requires a programming language (the codon table) and a set of programs written in that
68 language to enable it. How did physical chemistry enable discovery about 4Ga ago of the first genetic
69 coding rules from among so many others? Those rules were needed to write blueprints into the sequences
70 of the first genes whose translated products could then impose the same rules. The double helix ([Watson
71 and Crick 1953](#)) inspired a series of brilliant studies. These quickly revealed details of the genetic coding
72 table and the key roles of aminoacyl-tRNA synthetase (AARS)•tRNA cognate pairs ([Hoagland, et al. 1956](#);
73 [Berg and Ofengand 1958](#); [Nirenberg and Matthaei 1961](#); [Trupin, et al. 1965](#); [Jones and Nirenberg 1966a](#);
74 [Jones and Nirenberg 1966b](#); [Giegé 1972](#); [Ostrem and Berg 1974](#)).

75 Such rapid successes obscured questions about how the translation machinery might have arisen from the
76 prebiotic world. For decades, answers to such work remained as disconnected theories ([Crick 1968](#); [Yarus,
77 et al. 2009](#)). Most ([Koonin and Novozhilov 2017](#)), ignored the crucial role played by aminoacyl-tRNA
78 synthetases. As recently as 2000, and faced with very complex indications from the most recent
79 phylogenetic analyses, an authoritative review by experts including specialists in the area of AARS began
80 their conclusion as follows: “*It is unlikely that the aminoacyl-tRNA synthetases played any specific role in
81 the evolution of the genetic code; their evolutions did not shape the codon assignments.*” ([Woese, et al.
82 2000](#)). Our view is just the opposite. Their statement is logically the same as if we were to proclaim that
83 tRNA played no specific role in the evolution of the genetic code. The present-day code relies solely on the
84 specific fidelity of amino acid binding by AARS•tRNA cognate pairs. Neither statement accounts for the
85 way the present-day code relies solely on the specific fidelity of amino acid binding by AARS•tRNA
86 cognate pairs. Both imply that it emerged from a precursory system that vanished without leaving a trace.
87 That, in turn, would distinguish coding from virtually all other aspects of biology.

88 We chose to address the problem by excerpting model systems from full-length AARS. Those excerpts
89 include urzymes (~120-130 residues) and protozymes (~50 residues). They retain sufficient catalytic
90 proficiencies to enable us still to use enzymology at the limits of what we could retrieve about the past from
91 AARS structural biology ([Pham, et al. 2007](#); [Pham, et al. 2010](#); [Li, et al. 2011](#); [Carter 2014](#); [Martinez-
92 Rodriguez, et al. 2015](#); [Carter 2016](#); [Carter and Wills 2019b](#); [Carter 2022](#)). Only the most surprising aspect
93 of that work has been replicated ([Onodera, et al. 2021](#)). Nevertheless, we argue that such models for
94 ancestral AARS allow us to venture further into the past, beyond barriers that seemed evident to earlier
95 authors.

96 Translation requires a unique sequence of events. The first activates amino acids with ATP so that they can
97 join together spontaneously. The second links the amino acid covalently to tRNA. Doing so assigns a set of
98 chemical symbols—the anticodons—to represent each side chain. How did both activities arise at the same
99 time? That question is closely tied to a related one. The readout from genes to catalytic function is done by
100 proteins called aminoacyl-tRNA synthetases (AARS). What taught the ancestral AARS genes to enforce
101 those assignment rules? Indeed, the AARS gene products must read blueprints templated in the nucleotide
102 sequences of their own genes. That AARS can implement the rules first used to assemble themselves is a
103 property we call *reflexivity* ([Carter and Wills 2018b](#)). We will understand the origins of reflexivity only
104 when we can describe the earliest AARS•tRNA cognate pairs and the structural codes they used to recognize
105 both amino acid and RNA substrates ([Carter and Wills 2021b](#)). Their rates and specificities must also
106 support a cogent narrative that accounts for selection.

107 Recently, we have been able to form collaborative teams to improve the early model systems and leverage
108 them into a story that matches structural changes to increasing function ([Carter and Wills 2018b](#); [Carter
109 and Wills 2018a](#); [Wills and Carter 2018](#); [Carter and Wills 2019a](#); [Carter and Wills 2019b](#); [Wills and Carter
110 2020](#); [Carter and Wills 2021a, b](#); [Tang, et al. 2023](#); [Carter 2024](#); [Patra, et al. 2024](#); [Tang, et al. 2024](#)). At the
111 same time, we are enhancing those studies with new phylogenetics algorithms that sharpen the synthetase
112 and tRNA family trees ([Douglas, Bouckaert, Carter and Wills 2024](#); [Douglas, Bouckaert, Harris, et al. 2024](#);
113 [Douglas, Cui, et al. 2024](#)). From these and other studies we sketch an answer to the rhetorical questions

114 posed above. Nature had to discover the coding rules by teaching a set of genes how to read their own
115 blueprints. We summarize here a new paradigm that suggests we can now begin to understand how that
116 happened.

117 That paradigm includes the following elements.

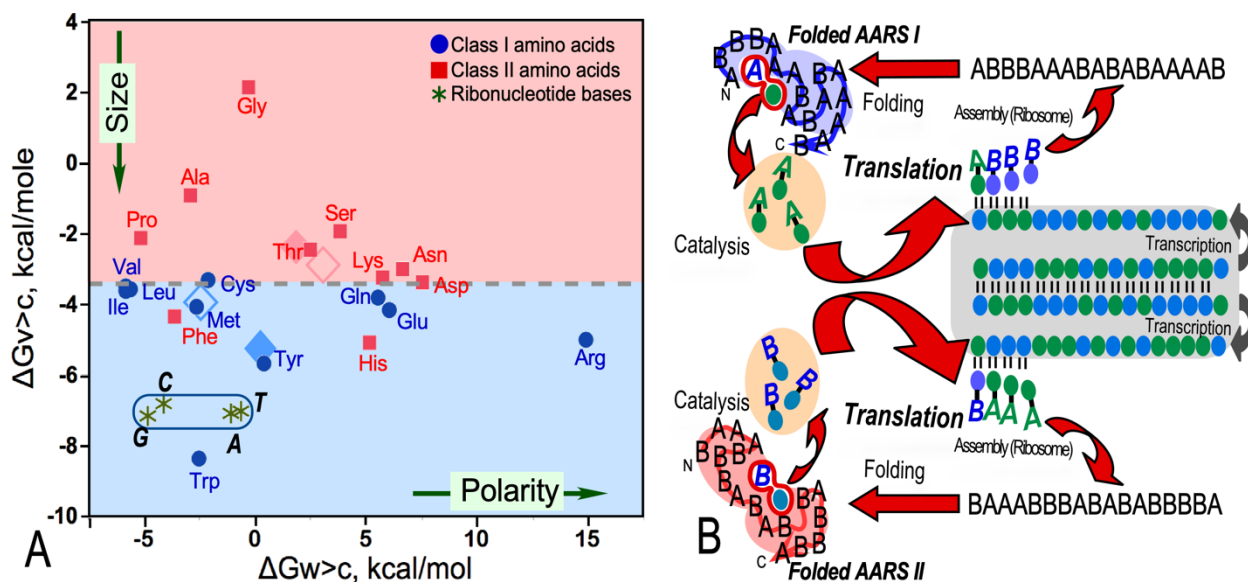
- 118 (i) The AARS are the central assignment catalysts necessary, along with their cognate tRNA
119 molecules, to link amino acids specifically to their symbolic representations in all protein-coding
120 genes. The earliest synthetase genes likely encoded Class I and II synthetases on opposite strands
121 ([Rodin and Ohno 1995](#)).
- 122 (ii) The genetic code assigns nonpolar and polar amino acids to complementary codons. That gives
123 protein sequences from bidirectional AARS genes inverted side-chain polarity ([Zull and Smith
124 1990](#); [Carter 2024](#)). Thus, the encoded Class I and II AARS pairs folded inside-out relative to one
125 another into different 3D structures. That, in turn, leads to different amino acid and RNA substrate
126 binding modes ([Carter and Wills 2019a](#); [Carter and Wills 2019b](#)). These data thus argue that initial
127 substrate recognition was rooted in the base pairing between their coding sequences. We reinforce
128 here the explanatory power of this observation ([Carter and Wills 2021b](#)).
- 129 (iii) AARS urzymes actually prefer to acylate T Ψ C-minihelices, instead of full-length tRNAs ([Tang,
130 et al. 2024](#)). Both protein and RNA components therefore appear to have functioned first as
131 simpler, single domains. Class I and II urzymes retain amino acid specificities consistent with
132 enforcing a code of ~4 letters. Studies of RNA recognition confirm details of the code used by
133 urzymes to recognize cognate minihelices ([Carter and Wills 2018a, 2019a](#)).
- 134 (iv) Catalysis by urzymes does not need active-site amino acids that were not present at early times
135 ([Tang, et al. 2023](#)). Phylogenetics and AI helped identify new AARS “urzymes” from novel
136 sources ([Patra, et al. 2024](#)). *E. coli* makes the same nested hierarchy of synthetase protozymes and
137 urzymes *in vivo* that it took two decades of analytical work to deconstruct.
- 138 (v) To be successful, further work requires new phylogenetic algorithms to solve unique problems
139 posed by the AARS family trees ([Carter, et al. 2022](#); [Douglas, Bouckaert, Carter and Wills 2024](#);
140 [Douglas, Cui, et al. 2024](#)).
- 141 (vi) Logic gating shaped the energy landscape on which coding emerged ([Carter and Wills 2021a](#)).

142 These aspects of our ongoing work now make a full platform for exploring the origins of genetics.

143 **Results**

144 **Synthetase genes must interpret their own blueprints.**

145 AARS acquired an enchanting natural property that was novel and apparently unique in the universe so far
146 known to us. In a manner akin to the training of artificial neural networks, they learned to read the
147 instructions for their own assembly. That property is called reflexivity (Fig. 1) because it entails self-
148 reference ([Hofstadter 1979, 2007](#)). Genetic coding stores the properties of the amino acid side chains—
149 their size and polarity—by matching them to a set of symbols (codons) (Fig. 1A). To implement coding,
150 Nature then put that physical chemistry to use. It discovered a set of gene sequences that initiates the
151 reaction cycle shown by the broad red arrows in Fig. 1B. Peptides with those sequences could fold in water
152 into 3D structures. Cavities in those 3D structures could recognize specific sets of amino acids and their
153 corresponding symbols. They also could speed up forming a chemical bond between them. Remarkably, at
154 the same time Nature also found how to immortalize those sequences as the ancestral AARS genes.



155

156 Figure 1. The origins of AARS reflexivity. **A.** Free energies of transfer for amino acids and ribonucleotide bases. These are the
 157 building blocks for proteins and nucleic acids. The Y axis is the free energy for transferring the side chain from vapor to
 158 cyclohexane. It is thus a surrogate for size. The X axis is the corresponding free energy for transferring the side chain from water
 159 to cyclohexane. It is a surrogate for polarity. The plot thus compares the physical chemistry of the nucleic acid and protein alphabets.
 160 Class I amino acids are blue dots; Class II amino acids are red squares. The colored background shows that Class I amino acids are
 161 predominantly bigger. They also span a larger range of polarity, although most are nonpolar Mean (solid) and median (outline)
 162 values for each Class are shown as diamonds of the same color. **B.** Schematic of the role of AARS in the information flow in
 163 genetics. A bidirectional ancestral gene and its mRNA transcripts are inside the gray panel. The gene is a bidirectional gene
 164 encoding Class I and II AARS on opposite strands. The respective translated peptides are written with a binary alphabet, with
 165 amino acids A, B activated respectively by Class I and Class II synthetases. Acylated RNAs are shown as capital letters linked to
 166 a green or blue ellipse, representing the symbolic codon representation. A folded conformation is essential for recognition of both
 167 amino acid and RNA substrates, and for stabilizing the two transition states for carboxyl group activation and acyl-transfer. Paired
 168 cycles of large red arrows define the reflexivity of AARS within each Class. Supplies of building blocks (acylated RNAs within
 169 amber ellipses) must be created by the two proteins, which must fold to catalyze the crucial reactions. Selection and gene replication
 170 are implicit in the cycle labeled “transcription”.

171 Converting the information in genes into proteins vastly expands it. The standard nucleic acid alphabet has
 172 only four letters, and these four letters have almost the same physical properties (Fig. 1A). Their sidechain
 173 volumes are both similar and larger than those of all amino acids except tryptophan. Their size greatly
 174 reduces the number of ways they can pack into tertiary structures. They also have almost the same polarity.
 175 Because both scales in Fig 1A are both logarithmic, the 20 (smaller) amino acids span a vastly greater range
 176 of both size and polarity. Thus, the structural and chemical roles of amino acids are far more diverse than
 177 those of the four bases. We have estimated on that basis that chemical engineering by proteins is roughly a
 178 billion times more diverse and proficient than that by ribozymes ([Carter and Wills 2018b](#); see
 179 [Supplementary material](#)). For that reason, the invention of coding signals a major advance in life’s agency.

180 AARS come in two different evolutionary Classes ([Eriani, et al. 1990](#)). Their respective specificities are
 181 shown in Fig. 1A as blue (Class I) and red (Class II). That partition has long puzzled the field. Many have
 182 tried to rationalize or explain the partition ([Delarue and Moras 1992](#); [Cusack 1994](#); [Rodin and Ohno 1995](#);
 183 [Ribas de Pouplana and Schimmel 2001a](#); [Ribas and Schimmel 2001](#); [Klipcan and Safro 2004](#); [Rodin and](#)
 184 [Rodin 2006](#); [Delarue 2007](#); [Safro and Klipcan 2013](#); [Takénaka and Moras 2020](#)). We consider the class
 185 partition to be an essential and penetrating clue to how genetic coding emerged, as outlined in the following
 186 section.

187 **The earliest synthetase genes encoded Class I and II synthetases on opposite strands.**

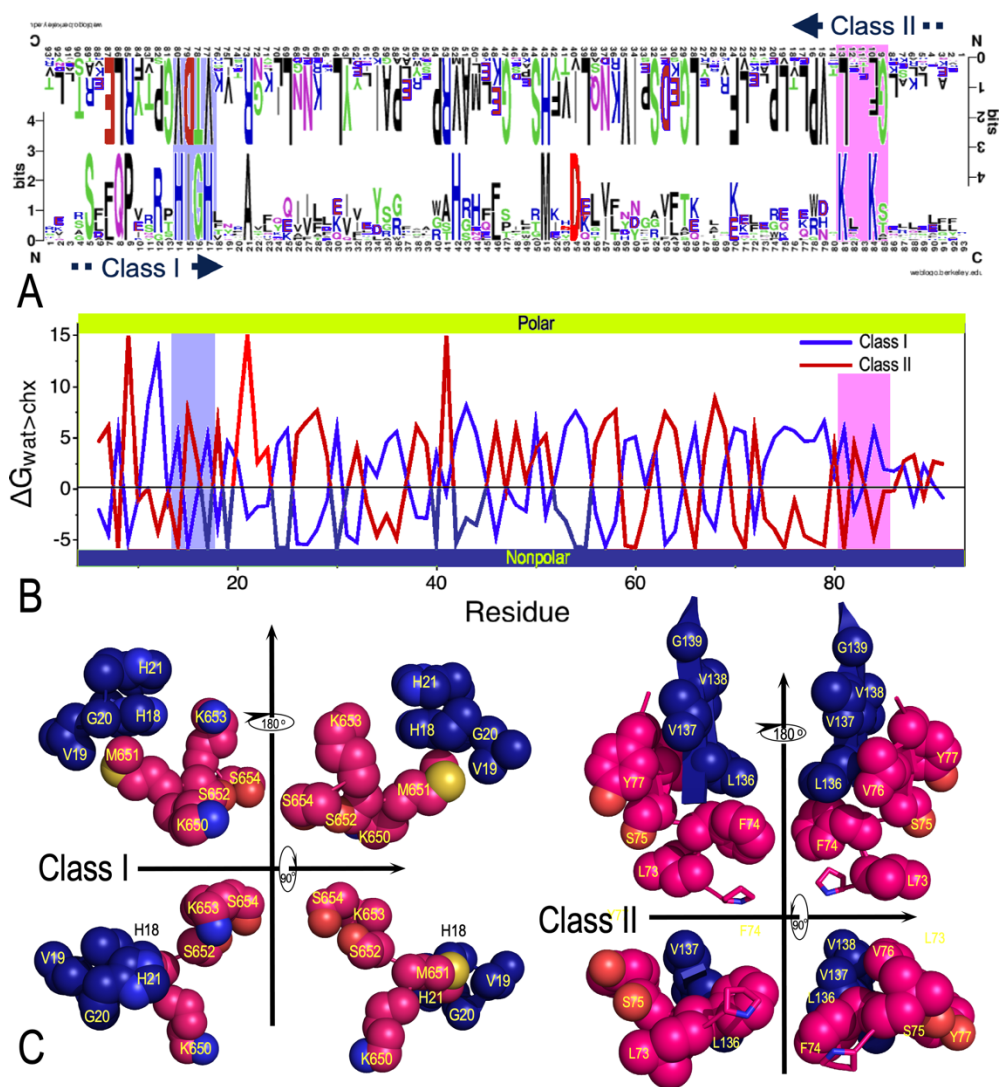
188 Rodin and Ohno proposed that the first Class I and II genes arose on opposite strands of the same nucleic
 189 acid ([Rodin and Ohno 1995](#)). They saw that antiparallel coding sequences for conserved “signature

190 sequences” defining the two Classes have far more base pairing than expected for random alignments. Three
191 orthogonal kinds of evidence validate the hypothesis: (i) The complementary regions do indeed encode
192 catalytically active synthetases AARS ([Carter, et al. 2014](#); [Carter 2015](#)). (ii) Independently reconstructed
193 ancestral Class I and II sequences show increasing base-pairing frequencies as they approach the root of
194 the tree of life ([Chandrasekaran, et al. 2013](#)). (iii) Bidirectional genes have been designed and expressed.
195 Translated products from both strands have the same catalytic proficiency, within experimental error, as
196 the native sequences ([Martinez-Rodriguez, et al. 2015](#); [Onodera, et al. 2021](#)). Finally, the Rodin-Ohno
197 hypothesis affords a potentially essential metric for constraining the Class I and II evolutionary trees. We
198 show below that the Rodin-Ohno hypothesis has extraordinary retrodictive power.

199 Two facets of bidirectional coding merit more comment. First, it suggests a contradiction: only one strand
200 of a gene carries unique sequence information. The inversion symmetry of base-pairing means that either
201 strand can be reconstructed if you know the sequence of the other. The prevailing view ([Crick 1970](#)) holds
202 for that reason that only one strand does the coding; the other is simply a template for reconstructing the
203 coding strand. How, then, can a bidirectional gene encode two radically different 3D protein structures? We
204 explore that question in the next section. Second, unlike all but one ([Wong, et al. 2016](#); [Takénaka and Moras
205 2020](#)) of the prevailing ([Crick 1976](#); [Carter 2008](#); [Yarus, et al. 2009](#); [Koonin and Novozhilov 2017](#)) ideas
206 about origins of coding, the Rodin-Ohno hypothesis serves us metaphorically, as an attractor ([Newman
207 1996](#)). It keeps spawning new, testable predictions. Each corroboration enriches the narrative.

208 **A curious inversion symmetry in the genetic coding table projects base-pairing into the proteome.**

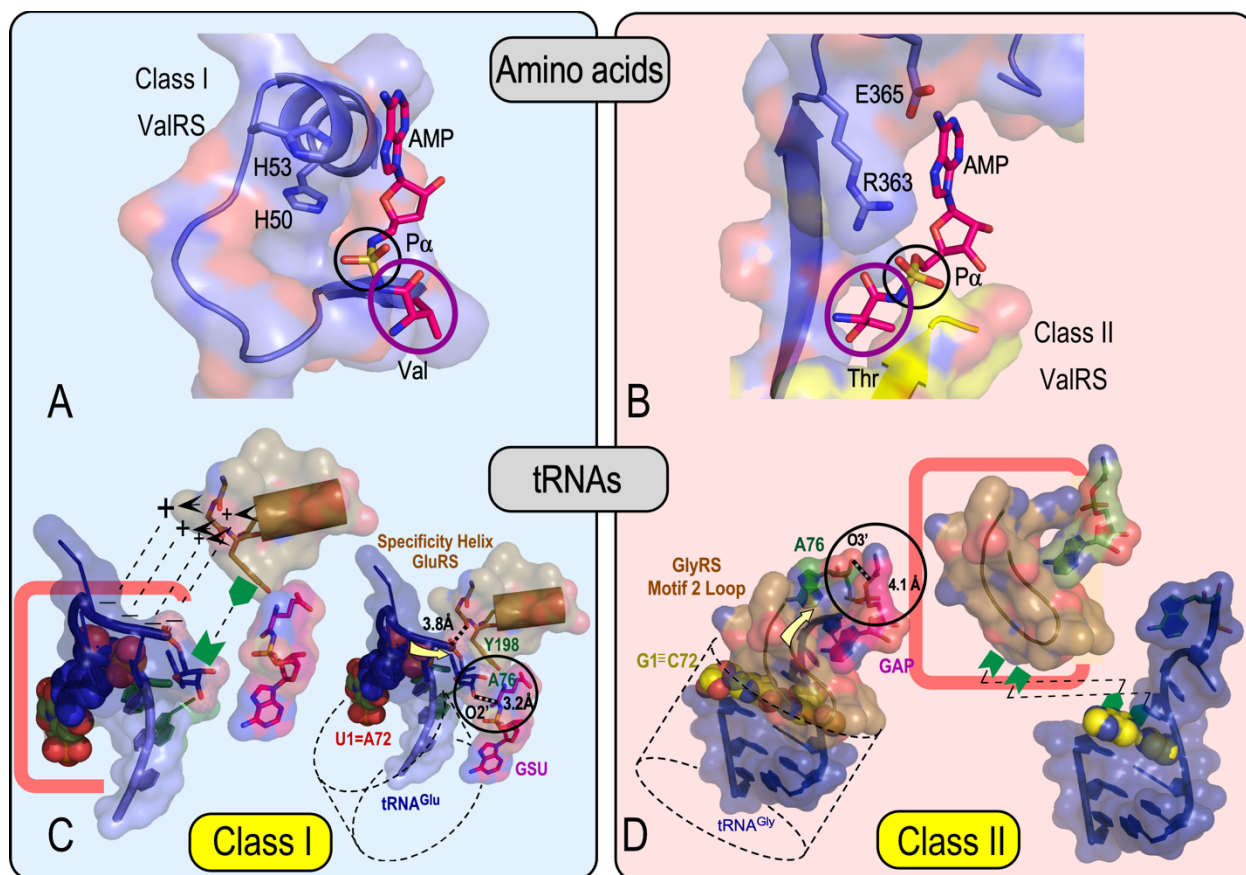
209 Zull and Smith ([Zull and Smith 1990](#)) observed that codon assignments are consistent with projection of
210 the nucleic acid inversion symmetry into the proteome (Fig. 2). Genetic codons are assigned to amino acids
211 such that all 14 codons for large, nonpolar residues are anticodons for highly polar residues. Dipeptides that
212 define turns (i.e. Pro-Gly; Gly-Pro) work in both directions on opposite strands. Sequence logos for a set
213 of bidirectional urzyme genes (Fig. 2A) illustrate the variety of sequences enabled by that detail of the
214 coding table. The mean residue by residue transfer free energy for water to cyclohexane partition of amino
215 acid side chains for that alignment, $\langle \Delta G_{w \rightarrow chx} \rangle$, reveals a stunning reflection symmetry in antiparallel plots
216 for the 81 amino acids in the designed bidirectional urzyme gene alignments (Fig. 2B). One consequence
217 of that reflection symmetry is that the two signature sequences in each Class (Fig. 2C) behave differently
218 in the folded conformations of Class I and II AARS. The Class I signatures, here HV(I)GH and KMSKS,
219 face into the active-site pocket where they accelerate activation and acylation by factors of 10^5 and 10^3 ,
220 respectively for the full-length enzymes. The two signatures are linked together only weakly by a nonpolar
221 interaction between the one hydrophobic side chain in each signature. In the absence of the two missing
222 domains, the catalytic contribution of the active site histidine and lysine residues is nearly absent. Their
223 catalytic roles require coupling imposed by domain motions ([Tang, et al. 2023](#)). By contrast, the Class II
224 signatures Motifs 1 and 2 form a compact nonpolar packing cluster that provides substantial structural
225 stabilization. In this sense, the ancestral Class I and II synthetase genes folded up inside out.



226

227 Figure 2. Inversion symmetry of base-pairing projects into the proteome. **A.** Sequence logos (Crooks, et al. 2004) from a set of
 228 designed bidirectional genes constrained by the active LeuAC “Goldilocks” and HisCA urzymes both of which have 93 residues.
 229 **B.** Plots of $\langle \Delta G_{\text{wat} \rightarrow \text{chx}} \rangle$, the free energy of transfer from water to cyclohexane for each sidechain in the sequence, an experimental
 230 metric for the polarity of the sidechain. Sequences of Class I AARS (blue) are plotted; those for Class II (red) are plotted right-to-
 231 left. **C.** Consensus Class-defining signature sequences from LeuAC (1WZ2) and HisCA show the impact of the inversion symmetry
 232 in the translated proteins. Class I signatures, H(I)GH and K(M)SKS, include mostly polar residues with key catalytic functions.
 233 Class II signatures, Motifs 1 and 2, form extensive nonpolar interactions.

234 *Inside-out folding retrodicts Class I, II amino acid and RNA substrate differentiation.* Class I, II protozymes
 235 provide the AARS ATP binding determinants. When superimposed on their adenosine moieties (Fig. 3A,
 236 B), we see that their α -phosphates are prochiral, hence offer non-equivalent binding environments. The
 237 amino acids react from opposite sides. Class I side chains, which are always bigger, have more room to
 238 grow because they face away from the protein (Fig. 3A). Class II side chains are smaller; they have less
 239 room because they would grow into the protein (Fig. 3B).



240

241 Figure 3. Inside-out folding accounts for both amino acid specificity and tRNA groove recognition. **A, B.** Superposition of the
 242 adenosine moiety of Class I ValRS and Class II ThrRS protozymes reveals that amino acid activation sites are diastereoisomeric.
 243 The loci of activated amino acids of Valine (A) and Threonine (B) face in opposite directions, accounting for the observation that
 244 Class I side chains are uniformly larger than homologous Class II side chains. **C, D.** Class I and II AARS recognize RNA substrates
 245 from opposite grooves. The 3'-acceptor stem forms a tight RNA hairpin (red rectangle) in Class I RNA substrates. The N-terminus
 246 of the Class I specificity-determining helix forms extensive interactions with that hairpin as indicated by dashed lines between
 247 protein positively charged groups and negatively charged phosphate groups in the RNA and nonpolar interactions between a protein
 248 aromatic residue and the nonpolar face of the A76 ribose (matching green blocks). Class II RNA 3'-acceptor stems are extended,
 249 leaving the initial base pair of the stem to form a platform for recognition by the motif 2 protein hairpin (red rectangle).

250 *Inside-out folding assures that Class I, II synthetases bind to opposite tRNA grooves.* A similar asymmetry
 251 differentiates tRNA binding by Class I and II AARS. The 3' DCCA terminus of Class I tRNA acceptor
 252 stems forms a sharp hairpin. That RNA hairpin is recognized by the N-terminus of the specificity
 253 determining α -helix in Class I urzymes. Class II acceptor stems are extended, leaving the initial base-pair
 254 as a platform for the Motif 2 protein hairpin. The synthetases for aromatic subclasses Ic (Trp, Tyr) and IIc
 255 (Phe) are exceptions ([Ribas de Pouplana and Schimmel 2001b](#)).

256 The duality of amino acid recognition combines with that of the RNA recognition mechanism. The 2D
 257 coordinate system in Fig. 3 (A,C vs B,D) underlies the initial discriminations allowing ancestral AARS and
 258 their RNA substrates to enforce a rudimentary coding alphabet. The distinctions between the signature
 259 sequences in Fig. 2C and the specificities defined in Fig. 3 reinforce the idea that the base-pairing inversion
 260 that is fundamental to heredity also projects into the proteome to define substrate recognition by ancestral
 261 AARS. We discuss the alphabets that might have grown from that projection in more quantitative detail in
 262 a subsequent section.

263 **AARS urzymes and protozymes are more than simple analytical constructs.**

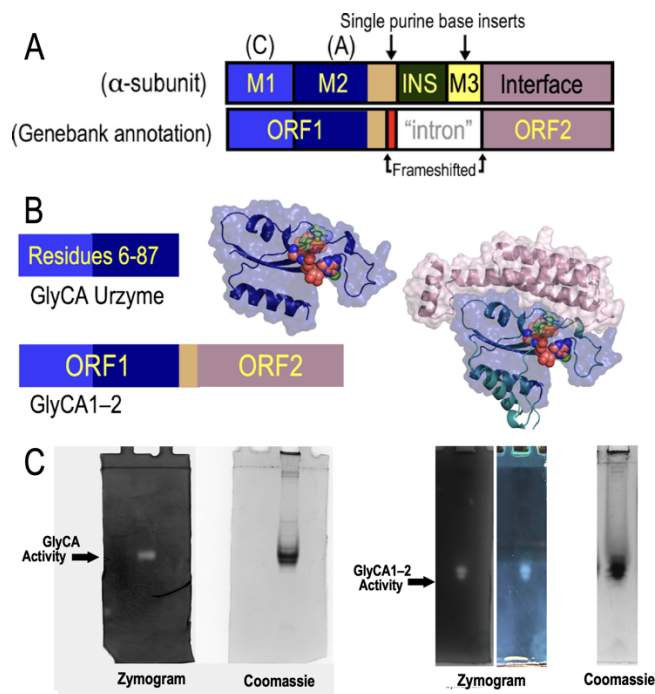
264 We coined the terms “urzyme” and “protozyme” to suggest that the constructs are models for ancestral
265 AARS forms. That notion implies that they may resemble evolutionary intermediates. We have argued that
266 their catalysis and substrate specificities are appropriate to have helped launch genetics. Despite an ample
267 literature, skepticism remains about their authenticity. Only one other group has tested any of our
268 biochemical work (Onodera, et al. 2021). In light of such lingering doubts, this section describes two recent
269 studies suggesting to us that the nested hierarchies may represent more than meets the eye. One entails an
270 urzyme for a new Class II AARS from an annotated eukaryotic genomic database (Patra, et al. 2024). The
271 other describes how *E. coli* generates a similar nested hierarchy of active segments.

272 *AARS urzymes occur in annotated genomic databases.* While assembling the aars.online database (Douglas,
273 Cui, et al. 2024), one of us (JD) noted that the Arctic Fox genome has a gene for the α -subunit from the
274 bacterium *Streptococcus alactolyticus* GlyRS-B. That curious orphan gene is a quirky result either of
275 horizontal gene transfer or contamination in the eukaryotic *V. lagopus* genomic data entry. The GlyRS-B
276 clade is found only in bacteria and chloroplasts. It is an unusual heterotetramer in which tRNA^{Gly} binding
277 is relegated to an idiosyncratic β -subunit with homology not to other AARS, but rather to four other cellular
278 proteins with RNA binding functions (Han, et al. 2023). Only the α -subunit is present in the genome. It is
279 not present in genomes from related foxes.

280 Two purine base insertions create a potentially functional intron in the annotated gene, separating it into
281 two ORFs (Fig. 4A). AlphaFold2 predicts that ORF1 has the same structure as the smallest Class II
282 synthetase urzyme derived from HisRS. A crystal structure for the *E. coli* GlyRS-B confirms the Alphafold
283 prediction (Han, et al. 2023). We expressed constructs with a shortened ORF1 and combining ORF1 with
284 ORF2 (Fig. 4B.) Both these excerpts accelerate aminoacylation of tRNA^{Gly} (albeit with high K_M values).
285 The GlyCA glycine activation rate acceleration is intermediate between those of Class II HisCA and Class
286 I TrpAC and LeuAC urzymes (Patra, et al. 2024).

287 Fig. 4C shows zymograms for both constructs with accompanying Coomassie-stained native PAGE gels.
288 These zymograms afford new, convincing visual evidence that the catalytic activity migrates with the major
289 band.

290

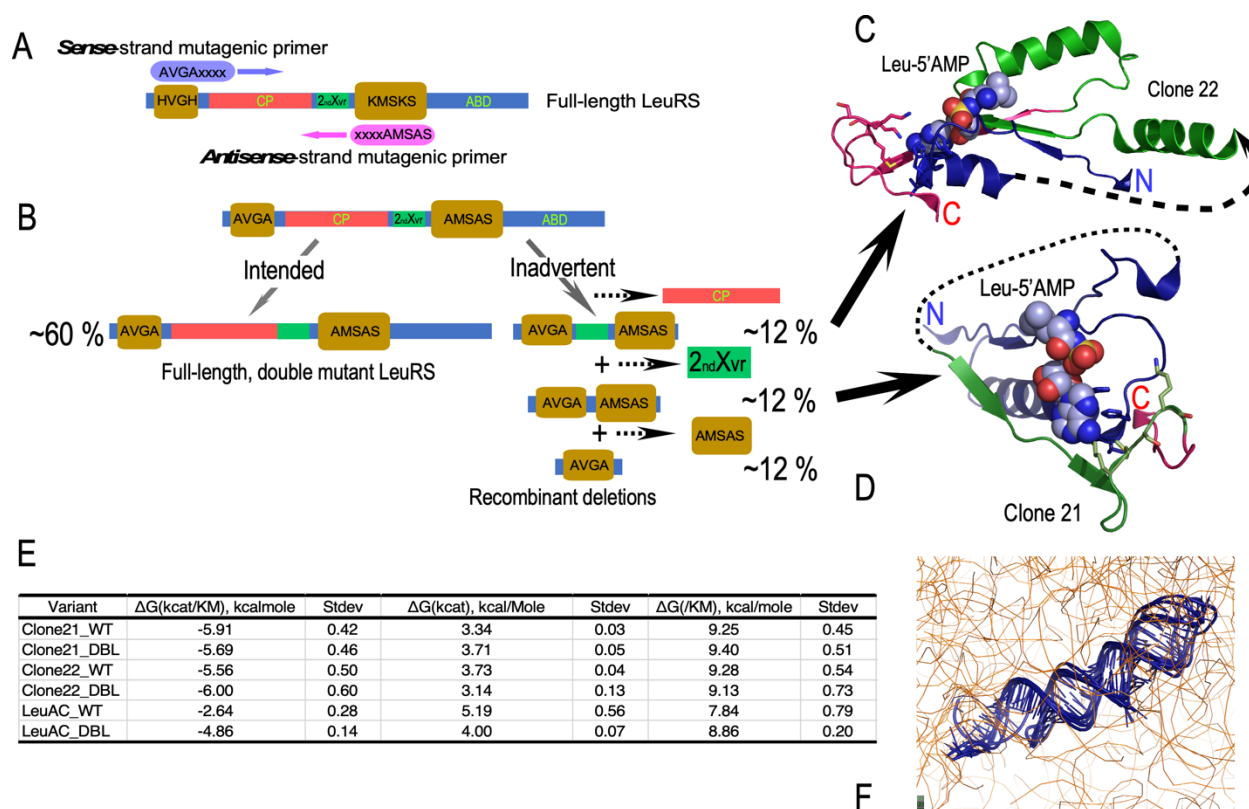


291

292 Figure 4. GlyRS urzymes derived from the *V. lagopus* genomic database (Adapted from (Patra, et al. 2024). A. Schematic of the
 293 annotated genomic entry for the Arctic Fox GlyRS-B α -subunit. Two inserted purine bases (top vertical arrows) create frameshifting
 294 (bottom vertical arrows) and an internal stop codon that produce two ORFs. ORF1 ends with a stop codon C-terminal to the red
 295 frameshifted sequence. B. These differences change the modularity of the *V. lagopus* GlyRS α -chain. 3D structures to the right are
 296 AlphaFold2 predictions for both ORFs. Colors are those used in the schematic. AlphaFold2 prediction matches closely that
 297 observed in PDB ID 7YSE (Han, et al. 2023). That allows visualization of likely binding geometry for glycine-5 sulfoamyl
 298 adenylate (spheres). The annotated intron coding sequence is frame-shifted and corresponds to what is called the insertion domain
 299 C-terminal to the active site. Motif 3 and an internal helix-turn-helix motif that covers the terminal adenosine in HisCA2 (Li, et al.
 300 2011) are thus missing in both GlyCA1 and Gly CA1-2. (Created using PYMOL (Pymol 2010)). C. Comparisons of a Coomassie
 301 stained native gel (right) to a zymogram visualized *in situ* using absorption of the MG orthophosphate complex by GlyCA and
 302 GlyCA1-2 (Cestari and Stuart 2013).

303 *E. coli* makes *in vivo* the same nested hierarchy of synthetase protozymes and urzymes the construction of
 304 which required two decades of analytical work. We designed the LeuAC urzyme by stitching together the
 305 conserved cores containing the HIGH and KMSKS signatures and deleting the connecting peptide and
 306 anticodon-binding domains. One of us (GQT) mutated the HVGH and KMSKS signatures to AVGA and
 307 AMSAS, thereby making a double mutant of full length native LeuRS for use in an earlier publication
 308 (Tang, et al. 2023). Only ~60% of the transformed *E. coli* colonies had the expected full-length double
 309 mutant plasmid (Fig. 5A). The remaining plasmids all had deletions of almost exactly those variable regions
 310 we had deleted *in vitro* (Fig. 5B). Urzymes bearing the two longer deletions catalyze both amino acid
 311 activation and aminoacylation. The rest of the recombinant deletions were protozymes. Both recombinants
 312 Clone 22 (Fig. 5C) and Clone 21 (Fig. 5D) retain the AMSAS sequence are active leucyl-tRNA synthetase
 313 urzymes (Fig. 5E).

314



315 Figure 5. *E. coli* provides an *in vivo* system to track modular evolution of AARS. A. Simplified scheme for mutating both HVGH
 316 and KMSKS signatures using bidirectional mutagenic primers. B. Schematic representation of plasmids sequenced from single
 317 colonies isolated after transforming *E. coli* with the full-length double mutant (AVGA/AMSAS) LeuRS plasmid. All recombinant
 318 deletions preserve the exact sequences from the input full length double mutant plasmid but are missing extended fragments from
 319 the LeuAC urzyme. C. Segments of Clone22, the “urzyme-like” deletion. It retains the second crossover connection (2nd Xvr) but
 320

321 is missing the C-terminal β -strand of the protozyme. It must fold differently from the urzyme because of the deletions, as indicated
322 by the dashed arrow. **D.** Segments of Clone21, the “Goldilocks” urzyme. Deletion of most of the specificity helix and second
323 crossover connection from the Rossmann dinucleotide-binding fold means that this construct folds differently because the green
324 segment containing the KMSKS loop is joined to the protozyme (dashed line). The 3D structures for **C** and **D** are taken from the
325 crystal structure 1WZ2 of the *Pyrococcus horikoshii* LeuRS. **E.** Steady state kinetic parameters of wild-type and double mutant
326 versions of Clone 21 and Clone 22 compared with those previously published for the corresponding LeuAC variants. **F.** Five 3D
327 structures predicted by Alfold3 ([Abramson, et al. 2024](#)) for the mRNA of full-length *P. horikoshii* LeuRS have been aligned by
328 superposing the coding region (dark blue) for the KMSKS-bearing loop (green segment in C). That region is the only motif shared
329 by all five predictions.

330 *The Goldilocks deletion maps to 3D structure in its mRNA.* The intermediate deletion is especially
331 interesting because although it retains the AMSAS sequence, it is missing most of the second crossover
332 connection of the Rossmann dinucleotide fold (dashed line in Fig. 5F). We already introduced the descriptor
333 “Goldilocks” for urzymes bearing this deletion because of its intermediate size. Its structure is unknown
334 but must differ substantially from that of the corresponding sequences in the full-length LeuRS (Fig 5D).
335 The Goldilocks deletion is an active, 81 residue synthetase. It represents the addition to the protozyme of
336 the mutant AMSAS signature.

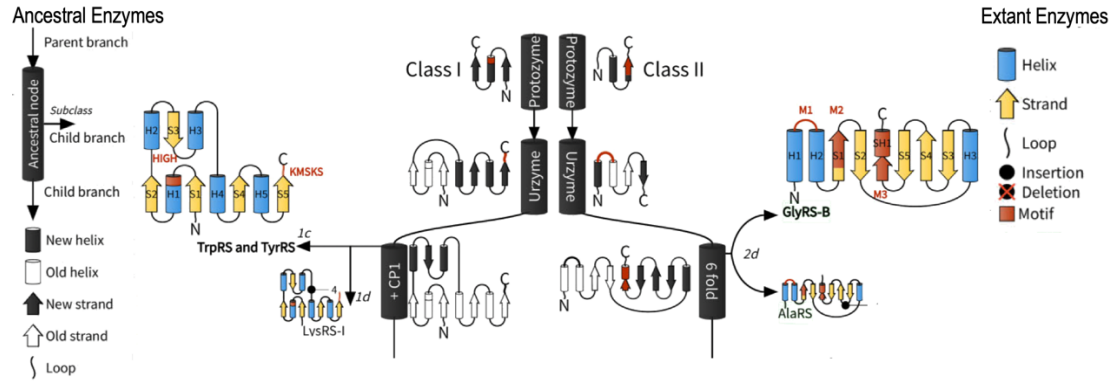
337 Notably, that short fragment seems to be both necessary and sufficient to enable catalysis of acyl transfer
338 from the activated aminoacyl-5’AMP to a substrate RNA. Further, as noted in Figs. 2A and 7, it is the same
339 length as the Class II HisCA urzyme. That parity appears also to resolve the challenge of designing Rodin-
340 Ohno bidirectional urzyme genes. For these reasons, the Goldilocks urzyme may also represent the minimal
341 evolutionary path from protozyme to urzyme.

342 The discrete locations and high reproducibility of these deletions strongly suggests that they result from 3D
343 structural properties of the gene itself. We had suggested as much in our initial paper describing the TrpRS
344 urzyme ([Pham, et al. 2007](#)). To investigate that possibility, we predicted mRNA structures using
345 AlphaFold3 ([Abramson, et al. 2024](#)). AlphaFold3 has only rudimentary capability to predict RNA structures
346 and the absence of many training datasets meant that predicted structures are highly variable. Nonetheless,
347 one long stem-loop was present in each prediction (Fig. 5F). That loop contains the exact codons for the
348 module that converts the LeuRS protozyme into the Goldilocks urzyme. That coincidence is consistent with
349 the hypothesis that the nested recombinant deletions in Fig. 5B, and by implication the synthetase
350 modularity arise from structures within the mRNA that promote recombination. *In vivo* models of this sort
351 could be an invaluable tool to study the enzymology of modular protein evolution (Fig. 5E).

352 **We need to adapt phylogenetics to address the unique problems posed by synthetases.**

353 Synthetases pose four problems for conventional phylogenetics algorithms: Each requires a qualitatively
354 different innovation in phylogenetic software. Our pursuit of more reliable ancestral reconstructions made
355 it necessary to implement solutions to each of these problems. We discuss these briefly in this section.

356 *AARS exhibit numerous, diverse insertions and deletions.* Indeed, both superfamilies likely grew to their
357 contemporary form from quite small protozymic ancestors largely by assimilating modular bits of genetic
358 information (Fig. 6). The gradual accumulation of insertion modules on the enzyme surface cannot be
359 treated adequately using amino acid substitution models alone ([Whelan and Goldman 2001b](#); [Le and](#)
360 [Gascuel 2008](#)). An additional Bayesian probability model is required to accommodate the birth and death
361 of insertion modules over long evolutionary timescales in both AARS superfamilies.



362
363 Figure 6. Insertion elements at the root of the AARS family trees [adapted from aars.online (Douglas, Cui, et al. 2024)]. Schematic
364 representation of two extant AARS (color) in the context of putative ancestors (black and white). The succession shown here is
365 truncated to include only the transition from protozyme to urzyme to catalytic domain. The purpose is to illustrate how important
366 modular accretion was to the evolution of both AARS superfamilies. Note the succession converting ancestral (black) to more
367 modern (white) secondary structures. Modular accretion continued to occur throughout the evolution of both superfamilies (see
368 Figure 4 of (Douglas, Bouckaert, Carter and Wills 2024)).

369 We addressed this problem in two stages. First, we created an online database with access to data on
370 secondary and tertiary structure alignments for both superfamilies (Douglas, Cui, et al. 2024). Many authors
371 in the field contributed to this database, which now provides more rigorously curated alignments. That
372 database helped us to identify how indels differ systematically from one synthetase to another. The second
373 stage was to build a Bayesian probability and a Dollo birth/death model for BEAST2 (Douglas, Bouckaert,
374 Carter and Wills 2024).

375 AARS for the most complex amino acids, tryptophan and tyrosine, resemble the oldest synthetases (Fig.
376 6). TrpRS and TyrRS appear to be the most recent additions to the Class I superfamily (Ribas de Pouplana,
377 et al. 1996). TrpRS and TyrRS also are the smallest of all AARS, and also have the simplest modularity
378 (Doublie, et al. 1995). Thus, they appear to be most ancient. Occam's razor suggests it is unlikely that earlier
379 versions had additional modules that all were lost in order to enhance their specificity. More likely, a
380 simpler dormant ancestral gene proved a better choice when amounts of the two amino acids increased to
381 the point where it became practical to create a new coding letter, as suggested in Fig. 6. We called such a
382 case "retrofunctionalization" to be compared with "neofunctionalization" and "parafunctionalization"
383 (Lynch and Force 2000; S. Rastogi and D. A. Liberles 2005). Other examples likely occur in both
384 superfamilies.

385 *AARS sequence evolution likely entailed substantial saltation.* One of the most problematic aspects of
386 building reliable trees is having to assume a gradualistic clock-like mutation process. Yet, branching often
387 appears to entail multiple correlated changes (Katsnelson, et al. 2019). Localized spikes of evolutionary
388 change have been well-documented experimentally for some more recent protein superfamilies (Bridgham,
389 et al. 2009; Manceau, et al. 2020). Such cases were summarized earlier at the macromolecular level
390 (Eldredge N 1972). Various attempts to solve this problem have had gradually improving success (Bokma
391 2002; Pagel, et al. 2006; Manceau, et al. 2020).

392 Consistent probabilistic treatment of rapid changes with birth/death model (Douglas, Bouckaert, Harris, et
393 al. 2024) led to several substantial improvements. We assumed that each lineage experiences a rapid
394 evolutionary spike, whose size is informed by the number of unobserved bifurcations along that lineage
395 (i.e., the number of stubs). Branching and evolution do indeed appear to be tightly coupled. Despite the
396 model's high dimensionality, several results convince us that it has many attractive properties. Most
397 important is enhanced statistical consistency, and a low covariation between rates and spikes. The model
398 consistently recovers correct known models from simulated data.

399 Unobserved speciation events do seem to have left their footprints on the lineages that have been observed.
400 We hypothesized that abrupt evolution could occur at any level in biology, demonstrating the process in

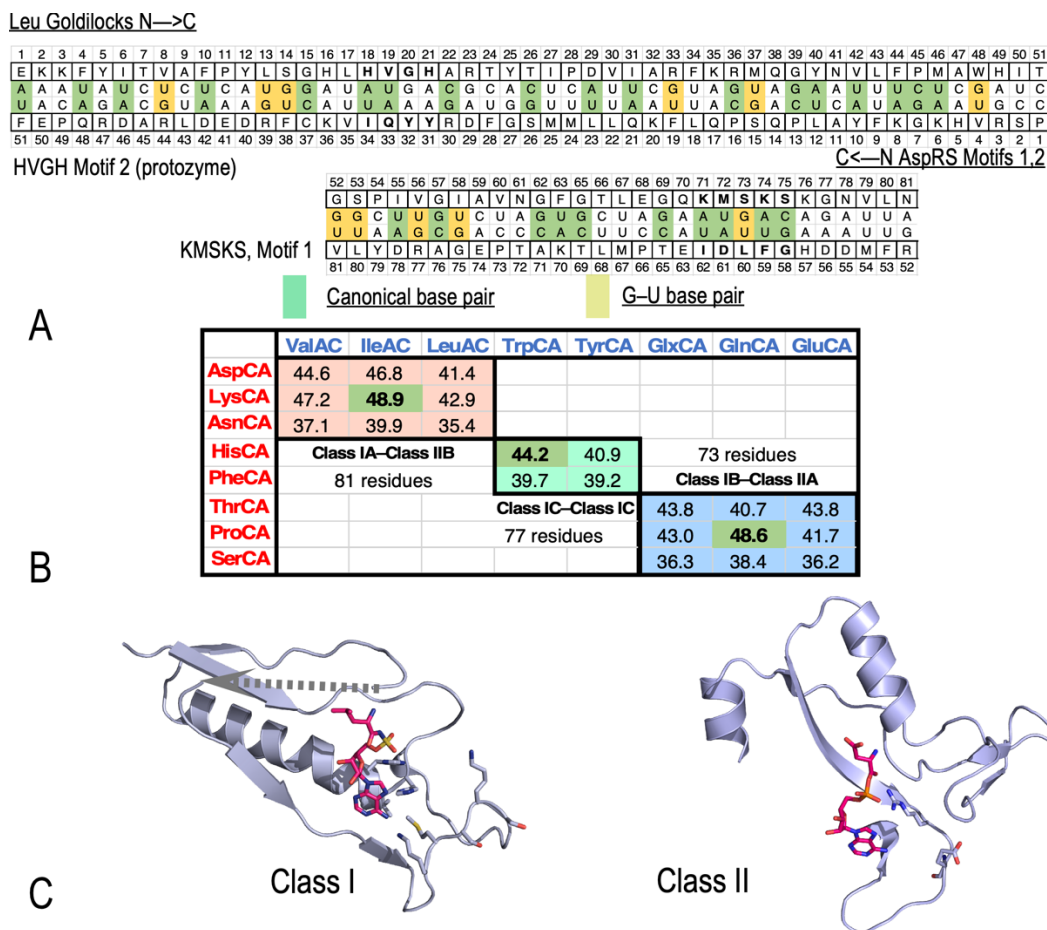
401 genes, morphologies, and languages ([Douglas, Bouckaert, Harris, et al. 2024](#)). Identifying them produced
402 more accurate divergence times. They also resulted in more reliable branchpoints. We demonstrated these
403 enhancements in three detailed case studies spanning the entire spectrum of phylogenies. Inferred AARS
404 molecular, morphological changes in cephalopods, and social evolution of Indo-European cognate word
405 phylogenies all changed differently. In particular, AARS branchpoints are indeed associated with highly
406 elevated fixation of point mutants.

407 *The amino acid substitution matrix must be updated dynamically.* Standard amino acid substitution models
408 rely on the critical assumption that the alphabet of amino acids remained constant through time ([Whelan
409 and Goldman 2001a](#); [Le and Gascuel 2008](#)). While this assumption might hold for younger proteins (e.g.,
410 among plants and animals), the assumption becomes increasingly unjustifiable as the age of the phylogeny
411 approaches the last universal common ancestor, and even more so dealing with the reflexive nature of the
412 AARS ([Carter and Wills 2018b](#); [Shore, et al. 2019](#)). This phylogenetic model, under active development,
413 involves the multiplication of a smaller, reduced $(n-1) \times (n-1)$ amino acid substitution matrix near the root
414 of the tree with the full $n \times n$ matrix in the present day, where n is the number of amino acids. Simulation
415 studies can often distinguish subfunctionalization (e.g., IleValRS \Rightarrow IleRS + ValRS) from
416 neofunctionalization (IleRS \Rightarrow IleRS + ValRS).

417 *Class I, II AARS trees must be coupled.*

418 Another basic issue for us is how to link the Class I and II AARS trees. We can make an overwhelming case
419 that even the earliest proteins needed both Class I and II AARS. The *sine qua non* of protein secondary
420 structure is binary patterns of either polarity and/or size ([Serrano, et al. 1992](#); [Muñoz and Serrano 1994](#)).
421 Polarity is even more essential for forming 3D structures ([Dill and MacCallum 2012](#); [Guseva, et al. 2017](#)).
422 The Rodin-Ohno hypothesis would, if verified, strongly constrain trees for the two superfamilies. We hoped
423 to infer such metric from the base-pairing frequencies of codon middle bases in the ~ 100 antiparallel
424 alignments of Class I vs Class II coding sequences ([Chandrasekaran, et al. 2013](#)). Codon middle-base
425 pairing between important segments of Class I and II urzyme genes could thus measure the strength of that
426 constraint. Assembling the data needed to compile those statistics, however, is hard owing to the numerous
427 subtle indels within and between different AARS families in both Classes.

428 After puzzling for many years, the Goldilocks LeuAC urzyme (Fig. 5) has told us that the 2nd crossover of
429 the Rossmann fold is not needed for aminoacylation. It was likely a later addition and may not have existed
430 during the brief stage of bidirectional coding. The Goldilocks recombinant provided a crucial datum by
431 revealing that only the short β -strand and KMSKS loop are required for aminoacylation. It has exactly the
432 same length as the recently characterized GlyCA Class II urzyme. That realization reduces our estimate of
433 how long the earliest ancestral AARS were to ~ 80 amino acids. Examples from work in progress are shown
434 in Fig. 7.



435
436 **Figure 7. Toward a Rodin-Ohno urzyme gene and a metric to link the early evolution of Class I and II AARS evolution.** **A.**
437 Antiparallel sequence alignments of designed 81-residue Class I and II Rodon-Ohno urzymes based on Class I *P. horikoshii* LeuRS
438 (1WZ2) and Class II *E. coli* AspRS (1C0A), with codon middle bases interleaved and colored according to base pairing, as
439 indicated. **B.** Partial table of codon middle-base pairing frequencies assembled from diagrams similar to that in **A.** Note in particular
440 the differences in length between the AspRS, GlyRS, and LeuRS urzymes and those for the other four. Also note that the Expected
441 pairing frequency between subclass IA LeuRS and subclass IIA GlyRS is substantially smaller than that with subclass II AspRS.
442 Thus, the subclass pairings suggested by middle base pairing frequencies are not as expected (see central part of the table). **C.**
443 Graphic images of the molecular structures implied by the alignments in **A.** Catalytic residues and both ammoacyladenylate ligands
444 are shown as sticks. The long, dashed arrow in Class I indicates a counterintuitive covalent bond connecting the C-terminal segment
445 to the protozyme.

446 This preliminary evidence replicates evidence from similar unpublished earlier work done before we found
447 the LeuAC Goldilocks urzyme. It suggests that three ancestral bidirectional genes originally paired Class
448 IA IleRS opposite Class IIB LysRS, Class II ProRS opposite Class I GlnRS, and HisRS opposite TrpRS. \

449 Completing the partial table in Fig. 7B may constrain the joint evolution of Class I and II AARS enough to
450 allow us to reconstruct ancestral sequences back to the time when urzymes enforced an earlier code of far
451 fewer than 20 amino acids.

452 These pairings are surprising because they suggest that wholesale sequence changes accompanied the
453 speciation as the earliest genes became increasingly specific AARS. Four of the putative ancestral
454 sequences shown in Fig. 7B resemble most closely those for “phase 2” amino acids—lysine, glutamine,
455 histidine, and tryptophan—that were not used in the earliest proteins {Wong, 2016 #833}. The
456 corresponding GlnAC LysAC, HisAC, and TrpAC urzyme sequences thus all seem to be examples of
457 retrofunctionalization ([Douglas, Bouckaert, Carter and Wills 2024](#)).

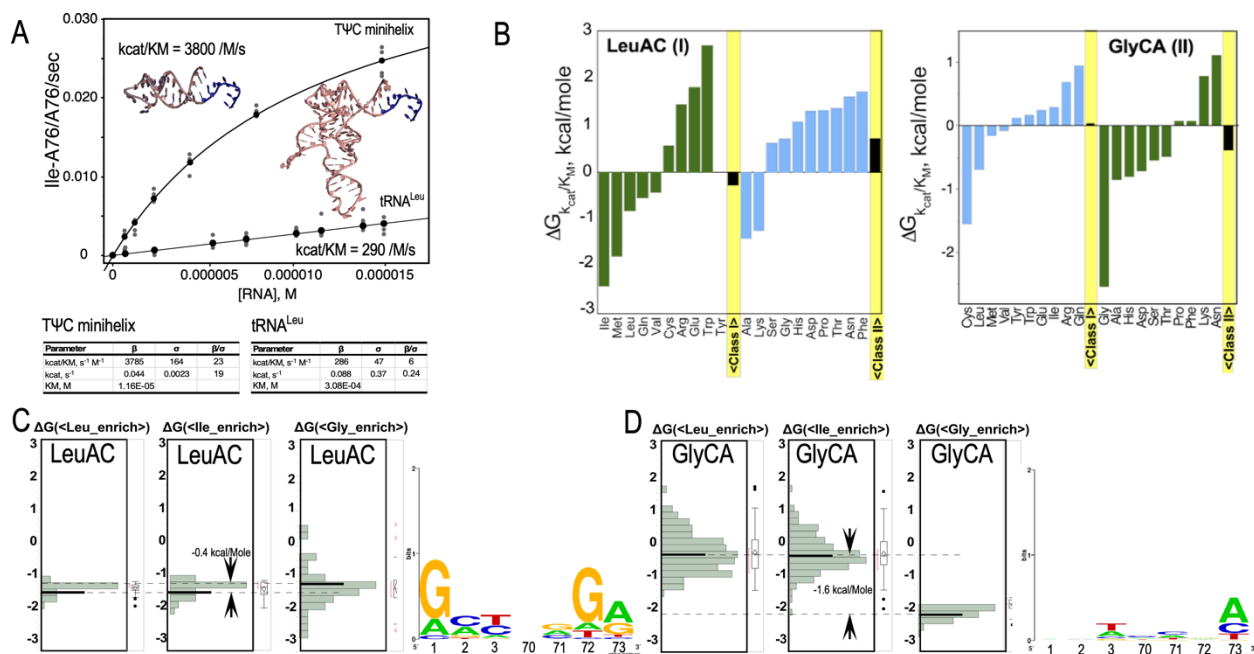
458 **AARS urzyme•minihelix cognate pairs are sufficiently functional to launch genetic coding.**

459 The phylogenetic studies outlined in the last section can open the door to more robust sequence distributions
 460 for AARS families. That is only part of the solution to our problem. Any account of survival will rely on
 461 relative fitness. The narrative we seek must also account for the survival of the succession of constructs.
 462 Proving fitness is a matter of enzymology. Rates and substrate recognition of each successive stage must
 463 be consistent with that narrative. This section recaps experiments that test how well AARS
 464 urzyme:minihelix cognate pairs might have functioned.

465 *Enzymological studies confirm that urzyme•minihelix pairs work better than either hybrid.* For decades, the
 466 tRNA anticodon was assumed to play the pivotal role in tRNA identity. It does, in ribosomal mapping to
 467 mRNA. However, a seminal proposal appeared in 1993 that today's dual-domain synthetase•tRNA cognate
 468 pairs grew from a much simpler set of single-domain•minihelix pairs ([Schimmel, et al. 1993](#)). Urzymes
 469 would not be discovered for another decade and a half ([Pham, et al. 2007](#); [Li, et al. 2011](#)). It took two and
 470 a half decades to work out details of the operational code. It took more than three decades to try to answer
 471 the question: can urzymes acylate not only tRNA but also minihelix substrates?

472 Schimmel's was a prescient idea, fueled by the fact that several groups had shown that full-length AARS
 473 could acylate minihelices that lacked the anticodon ([Francklyn and Schimmel 1989](#); [Schimmel 1991](#);
 474 [Frugier, et al. 1992](#)). The idea that the anticodon and corresponding anticodon-binding domains could have
 475 come later also implied that the original basis for tRNA recognition by ancestral synthetases must have
 476 depended on a different code embedded in the acceptor stem. The authors proposed no details of how such
 477 an operational RNA code might work. Yet, it was the skeleton of the crucial notion that the ancestral coding
 478 system might originally have worked with much simpler machines.

479 For any system to evolve, it has to function. Any ancestral system must have been able to perform two tasks
 480 with some degree of substrate specificity: activate amino acids and acylate RNA substrates. It is crucial to
 481 know whether or not urzymes can recognize amino acids and acylate tRNA minihelices. We recently
 482 answered this question dramatically. Urzymes not only have the capability to acylate cognate minihelices;
 483 they actually prefer the minihelix substrate by about tenfold (Fig. 8A).



484
 485 **Figure 8. Evidence for the evolutionary fitness of AARS urzyme•minihelix cognate pairs.** **A.** Comparison of Michaelis-Menten
 486 kinetic data for the LeuAC urzyme acylation of TΨC-minihelix^{Leu} and tRNA^{Leu} [adapted from ([Tang, et al. 2024](#))]. **B.** Amino acid
 487 specificity spectra for the LeuAC and GlyCA urzymes [adapted from ([Patra, et al. 2024](#))] **C.** Note in particular that although the

488 Class I and II Urzymes cannot distinguish amino acids as well as full-length AARS, they tend to favor amino acids from their own
489 Class. **C, D.** RNA specificities of LeuAC and GlyCA urzymes. Histograms for the enrichments of aminoacylated sequences over
490 those inferred from the total library for minihelices acylated with Ile and Leu (**C**; 39 hyperacylated sequences) by AVGA LeuAC
491 mutant and Gly (**D**; 204 hyperacylated sequences) by GlyCA. The vertical axes are expressed as free energy values in kcal/mole,
492 computed on the basis of the enrichment. Black horizontal lines denote the average values. Web Logos compare the RNA specificity
493 of the AVGA LeuAC mutant and GlyCA enriched fractions. The histograms and web logos represent 10-fold (**C**) and 22-fold
494 hyperacylated (**D**) fractions, respectively.

495 The Michaelis-Menten plot for LeuAC urzyme (Fig. 8A) shows that the Leucine minihelix is an order of
496 magnitude better substrate than tRNA^{Leu}. Neither hybrid protein•RNA system works as well. Thus, if the
497 RNA substrate is a minihelix, the best enzyme to acylate it is the cognate urzyme, and conversely. LeuAC
498 has a marked preference the minihelix and LeuRS much prefers tRNA^{Leu}. That antisymmetry is unexpected
499 evidence that the two modules co-evolved from the start of genetic coding.

500 Experimentally, the LeuAC and GlyCA urzymes also differentiate complementary subsets of ~5 similar
501 amino acids (Fig. 8B). The mean free energy difference between the correct and incorrect amino acid Class
502 is about 1 kcal/mole. That means the urzymes mistakenly activate an amino acid from the wrong class about
503 1 in 5 times. That specificity is comparable to what we observed previously ([Carter, et al. 2014](#)).

504 The same two urzymes also acylate different subsets from a minihelix library (Fig. 8C, D). That library
505 contained all combinations of the seven bases at the top of the acceptor stem. Subsets acylated by LeuAC
506 and GyCA were biotinylated and collected on streptavidin beads. The selected minihelices were then
507 sequenced, as was the un-acylated library itself. As noted in **Methods and Materials**, we used the same
508 library for each urzyme. That library had acceptor stem sequences derived from tRNA^{Leu}. For that reason,
509 the two urzymes behave differently. We show in Fig. 8C, D the distributions of the most hyperacylated
510 minihelices from the library, together with web logos for those distributions. The GlyCA urzyme acylated
511 fewer minihelices with greater enrichments. It appears somewhat more selective for minihelices from the
512 library, likely because the constant region of the acceptor stem has some residual specificity.

513 As with the amino acids, cross-specificity is evident. The mean differences in free energies are ~ -1
514 kcal/mole for activation by both urzymes for class-specific acylation and -0.4 kcal/mole and -1.6 kcal/mole
515 for minihelix specific aminoacylation by LeuAC and GlyCA.

516 This experimental platform enables us now to trace the growth of the coding table. As far as we know these
517 are the only data available so far on the likely fidelity of very early translation systems. They raise this new
518 question: Are these rates and fidelities sufficient to compose a self-consistent reflexive system? Answering
519 that question appears now to be within reach. The phylogenetic software set out in the previous section
520 should enable us to estimate which amino acids were in likely ancestral coding tables. They also should
521 allow us to construct reliable ancestral distributions of sequences for the relevant AARS that might have
522 enforced such ancestral codes. We can then construct combinatorial libraries for those very AARS. We
523 have begun to express and characterize such libraries (Patra, S.K. *et. al.* in preparation). We also have
524 written software for treating the kinetics of catalytic mixtures ([Douglas, Carter and Wills 2024](#)).

525 **Reciprocally coupled gating suggests the existence of two, new, biological forces.**

526 The reason there are so many different attempts to define the origin of life is that life itself seems so highly
527 improbable. For life to emerge on early post-Hadean earth many, quite different things all had to happen at
528 almost the same time. Translation was only one of the most important and challenging of those things. It
529 could have emerged only under quite special circumstances. Those include a steady input of free energy
530 ([Liu, et al. 2020](#)), rudimentary metabolism ([Sobotta, et al. 2020](#); [Stubbs, et al. 2020](#)), mutualism relating
531 polypeptides and nucleic acids ([Lanier, et al. 2017](#)), and possibly some sort of compartmentation ([Zhu, et
532 al. 2013](#)). These areas are outside the scope of this work, but they all must have been provided in some
533 form.

534 Genetic coding alone seems, at first glance, to be highly improbable. Afterall, within the vast combinatorial
535 space of amino acid-to-codon assignments and gene sequences, only a tiny fraction can enforce reflexive

536 genetic coding. With more than 10^{84} ways to assign 64 codons to one or more of the 20 amino acids and
537 $>10^{105}$ possible amino acid sequences that are the length of an urzyme, an exhaustive search for suitable
538 subsets of assignments and genes would exceed the resources of the known universe. How then did a
539 suitable combination emerge here on earth?

540 Such combinatorial problems abound almost anywhere one looks; they are especially common to biology.
541 They can be described as “Levinthal-like” paradoxes because, they share the formulation articulated by
542 Levinthal for protein folding ([Levinthal 1968](#); [Dill and Chan 1997](#)). Solutions to the Levinthal paradox
543 take the form of a free energy landscape that adopt the form of a “funnel” ([Socci, et al. 1998](#)). The idea of
544 a funnel is that the gradient of the free energy surface guides the search through combinatorial space by
545 making changes favorable if they move toward a free energy minimum. Analogously, the search for suitable
546 combinations of amino acid-to-codon assignments and genes may be resolved by a fitness landscape that
547 takes the form of a funnel that guides the search through combinatorial space towards combinations that
548 support a robust translation system.

549 Such combinatorial problems abound almost anywhere one looks; they are especially common to biology.
550 They can be described as “Levinthal-like” paradoxes because, they share the formulation articulated by
551 Levinthal for protein folding ([Levinthal 1968](#); [Dill and Chan 1997](#)). Solutions to that paradox take the
552 general form of a free energy landscape or “funnel” ([Socci, et al. 1998](#)). The idea of a funnel is that the
553 gradient of the free energy surface guides the search through combinatorial space by making changes
554 favorable if they move toward a free energy minimum. An analogous fitness landscape likely guided the
555 search for suitable combinations of amino acid-to-codon matching and genes.

556 What might have shaped such a funnel in the fitness landscape on which genetic coding emerged? A force
557 is a change in an energy field that changes the direction or velocity of motion. If we view evolution as
558 motion on a landscape then selective constraints imposed on that motion can be seen as biological “forces”.
559 One selective constraint is that a translated gene sequence can affect chemistry best if and only if it always
560 folds into about the same 3D structure. Another is that a sequence will fold if and only if its amino acids
561 occur in an appropriate order.

562 These two constraints are shown by the pairs of horizontal red arrows in Fig. 1B. Each acts as a logical
563 gate, reducing the passage either of peptides that do not fold, or of sequences that do not obey the code.
564 The two logical gates can exert an extraordinary, iterative force if the consequent of the first serves as the
565 antecedent of the second and *vice versa*. We argued earlier ([Carter and Wills 2018b](#); [Wills and Carter 2018](#))
566 that such coupling makes the launch of coding from protein much more probable than from ribozymal
567 synthetases.

568 Searching at the same time for the coding table and the gene sequences required to implement it also helps
569 solve the Levinthal paradox by shaping the landscape of evolution. Computer modeling of autocatalytic
570 sets shows that side reactions known as “parasites” are the chief threat to survival ([Takeuchi, et al. 2017](#)).
571 Reciprocally coupled gating eliminates parasites in both directions.

572 When two logic gates are coupled tip-to-tail with the consequent of one serving as the antecedent for the
573 other, they make a “strange loop” ([Hofstadter 1979, 2007](#)). Tip-to-tail coupling creates a self-referential
574 cycle. Self-reference establishes a threshold beyond which consistent systems become incomplete.
575 Incompleteness, in turn, implies a reservoir of possibilities. The strange loop thus creates a separate force.
576 In chemistry the spatial gradient of a species’ chemical potential induces a change in its equilibrium
577 distribution. That gradient of novelty drives discovery in ways analogous to the role of chemical potential.

578 Searching at the same time for code words (the coding table) and code keys (the AARS gene sequences)
579 adds a cooperative element to the search that enhances its efficiency ([Carter and Wills 2018b](#); [Wills and
580 Carter 2018](#)). At the same time, the implicit incompleteness also ensures the emergence of novelty. Novelty,
581 in turn, is the grist on which selection and evolution work.

582 **Conclusions**

583 How do the new results in this paper help answer the questions posed in our opening paragraph? We crafted
584 them with some care, because the authors differ on how to view the question of Nature's agency. When it
585 comes to evolution, agency offers a helpful shorthand to keep sentences from being too verbose. Too much
586 agency can be conflated with creationism. In the central chapter of his book *How Life Works: A User's*
587 *Guide to the new Biology*, Philip Ball outlines why this question is so vexing (Ball 2023). Here, in a similar
588 vein, we try to emphasize that the question is also intrinsic to our subject. The birth of genetic coding was
589 indeed a remarkable part of biology's transcendence of chemistry. As such, we must recognize how close
590 to optimal the genetic system really is.

591 Nor can we ignore what appears to be a sense of purpose. Lineages of ever better enzymes built, expanded
592 and refined a table of symbols to capture the physical chemistry of amino acids. As the table grew, so too
593 did the capacity to enhance the catalytic activity and specificity of the genes whose translated products
594 enforced that growing table. Some property of the first genes allowed that to happen.

595 A growing set of new pilot studies help tie the early evolution of AARS to the emergence of genetic coding.
596 *E. coli* can generate *in vivo* a nested set of active excerpts similar to the protozymes and urzymes we had
597 previously designed as models for ancestral AARS. A mid-sized "Goldilocks" variant made *in vivo* and a
598 short new Class II urzyme suggest how to design bidirectional urzyme genes. The inversion symmetry of
599 base pairing in sense/antisense ancestral AARS genes projects into the proteome. That leads in turn to
600 contrasting amino acid and RNA substrate binding modes. Codon middle-base pairing can help coordinate
601 the building of Class I and II AARS family trees. Acylation of specific subsets of a minihelix library by
602 Class I and II AARS urzymes shows RNA substrate specificity for the first time. Finally, new phylogenetics
603 routines solve four problems that had blocked rooting AARS trees in reduced coding alphabets. Those
604 algorithms will enhance our experimental work. Together, they form a viable platform to study how Nature
605 likely built the earliest genetic coding tables and began to enforce them with quite simple AARS
606 urzyme•minihelix cognate pairs.

607 **Methods and Materials**

608 **Zymography**

609 To visualize amino acid activation in a native polyacrylamide gel chromogenically, Zymography was done
610 using a 1.5 mm thick native gel of 8% resolving and 5% stacking which devoid of SDS. Protein samples
611 were prepared without adding SDS and β -ME to maintain the native conformation of protein. Prepared
612 protein (50 μ g) was then loaded into two separate wells of native gel. Electrophoresis was done with 40 mA
613 steady current at 4°C. After electrophoresis for 30 minutes the gel was transferred from the glass plates for
614 staining. Then the gel was shaken twice with double distilled water for 5 min. After washing, reaction
615 buffer—50 mM HEPES pH 7.5, 100 mM amino acid (here glycine), 20 mM MgCl₂, 50 mM KCl,
616 pyrophosphatase solution [NEB] (0.1 Unit/ml) and polyethylene glycol (PEG-8000; Sigma-Aldrich, Cat.
617 No. 25322-68-3) to a final concentration of 5% (w/v)—was poured over the gel. The gel setup was shaken
618 for 45 minutes at 4°C. This step ensures complete soaking of substrate mixture into the gel; the low
619 temperature reduces inactivation of the enzyme during perfusion.

620 After the perfusion most of the solution was decanted, leaving minimal solution in a static condition at
621 37°C. Synthetase activity was activated by adding 5mM ATP solution dropwise onto the gel surface to cover
622 the whole gel. The gel was incubated for 30 minutes. After decanting the reaction mix from the gel box, the
623 staining solution (0.05% Malachite green in 0.1 N HCl and 5 % hexa-ammonium heptamolybdate
624 tetrahydrate solution in 4 N HCl) was added directly onto the gel box. Staining was promoted by shaking
625 for 2 minutes on a gyro-shaker. The staining solution was made as described in (Onodera, et al. 2021). The
626 gradual development of green bands (620 nm) around GlyCA protein present in the gel signified the
627 phosphomolybdate-malachite green complex formation and thus the *in-situ* activity of amino acid activation
628 by GlyCA. The gel was photographed using Gel Doc™ XR+ from BIO RAD imaging machine.

629 **Isolation of recombinant deletions from a plasmid containing a LeuRS double mutant.**

630 We observed recombinant deletions by sequencing plasmids from ~60 colonies of *E. coli* transformed with
631 a double mutant of WT *P. horikoshii* LeuRS. Histidine and lysine residues of both HVGH and KMSKS
632 catalytic signatures were mutated to alanine for an earlier project ([Tang, et al. 2023](#)).

633 *Identification of in vivo recombinant deletions.* The plasmid vector was pET11a (Novagen, Sacramento,
634 CA, United States). DNA oligos were ordered from IDT (Integrated DNA Technologies, Coralville, IA,
635 United States). Phusion™ Plus PCR Master Mix (Cat# F631S) was purchased from Thermo Fisher
636 Scientific (Waltham, MA, United States). *E. coli* competent cells were partially from Agilent (XL10-Gold
637 ultracompetent cells, Cat# 200315, Santa Clara, CA, United States), and partially using DH5a, home-made,
638 prepared according to the method described by Sharma in 2017. Restriction enzyme DpnI (Cat# 500402,
639 United States) was purchased from Agilent. Purifications and handling of DNA fragment and plasmid used
640 the QIAquick Gel Extraction Kit (Cat# 28704, Qiagen, Hilden, Germany) and QIAprep Spin Miniprep Kit
641 (Cat# 27104, Qiagen, Hilden, Germany) according to instruction, unless stated otherwise.

642 *Method for preparation of double mutant plasmid.* DNA oligo sequences were designed according to the
643 manual of QuikChange II-E Site-Directed Mutagenesis Kit (Santa Clara, CA, United States) as published
644 online at <https://www.agilent.com/cs/library/usermanuals/public/200555.pdf>. In line with the oligo DNA
645 design, and for creating the double mutant plasmid DNA, we have implemented a megaprimer PCR method
646 described by Picard ([Picard, et al. 1994](#)).

647 After DpnI digestion of the PCR mixture at 37°C from hours to overnight according to instruction provided
648 by Agilent (Santa Clara, CA, United States), the DpnI-digested PCR mixture is ready for *E. coli*.
649 transformation. The transformation was done using 0.3 mL DpnI-digested PCR mixture per transformation.

650 *E. coli. transformation and plasmid sequencing.* *E. coli* transformation was proceeded according to
651 conventional procedure ([Fritsch, et al. 1982](#)) and spread onto LB (Lysogeny broth) agar plates with
652 antibiotics (Ampicillin 50 mg/mL, or Carbenicillin 100 mg/mL). After incubation overnight at 37 °C,
653 individual colonies on the LB agar plate were readied for miniprep by following the standard protocol
654 described by QIAprep Spin Miniprep Kit (Cat# 27104, Qiagen, Hilden, Germany)

655 Plasmid DNA samples were sequenced following standard sequencing sample preparation and submission
656 described by Eton Bioscience (San Diego, CA) online:

657 [https://www.genewiz.com/en/Public/Resources/Sample-Submission-Guidelines/Sanger-Sequencing-](https://www.genewiz.com/en/Public/Resources/Sample-Submission-Guidelines/Sanger-Sequencing-Sample-Submission-Guidelines/Sample-Preparation#sanger-sequence)
658 [Sample-Submission-Guidelines/Sample-Preparation#sanger-sequence](https://www.genewiz.com/en/Public/Resources/Sample-Submission-Guidelines/Sample-Preparation#sanger-sequence) with minor modifications. Basically,
659 each sequencing sample was in 15 mL, containing 100-200 ng plasmid DNA with 5-10 picomole
660 sequencing primer, i.e. seqn1_AVGA or AMSAS_seqn2rv as detailed sequences in Table 1 above.

661 *Method for plasmid sequencing and sequencing analysis.* Sequence analysis was performed using
662 conventional methods as follows. First, all the DNA sequences were translated into 6 reading frames in
663 amino acid sequence using Expasy translate online at <https://web.expasy.org/translate/>. For screening the
664 tentative urzyme candidates, each amino acid sequence from 6 reading frames from a single sequenced
665 plasmid DNA sample were aligned with the sequence of urzyme template using conventional sequencing
666 alignment software such as MAFFT online at <https://www.ebi.ac.uk/jdispatcher/msa/mafft?type=protein>
667 ([Kato, et al. 2019](#)) or MUSCLE online at <https://www.ebi.ac.uk/jdispatcher/msa/muscle?type=protein>
668 ([Edgar 2004](#)). Novel plasmids were designated urzyme, urzyme-like and goldilocks, according to their
669 sequence content and were collected and subjected for further downstream identification at the enzymology
670 level after protein biochemistry procedures.

671 **Aminoacylation of minihelix combinatorial libraries**

672 We constructed a minihelix library with sequences 35 nucleotides long derived from the acceptor-stem and
673 TΨC stem loop of tRNA^{Leu}. We varied only bases 1-3, 70-73 and these were given all ~16.4 K possible

674 combinations. The acceptor stem of tRNA^{Leu} is very GC rich and so will bias recognition by different
675 synthetases.

676 Minihelix library constructs were generated (Integrated DNA Technologies) with seven fully randomized
677 positions in the acceptor stem including the first three base pairs and the Discriminator base ([Schimmel, et
678 al. 1993](#)). Minihelix RNA was heated to 90°C for 2 minutes and cooled to 22°C for 10 minutes. Refolded
679 minihelix RNA was resuspended in acylation buffer: 50 mM HEPES, pH 7.5, 20 mM MgCl₂, 20 mM KCl,
680 5 mM DTT, 20 mM ATP and 20mM amino acids. Aminoacylations were performed by mixing the minihelix
681 RNA with a specified Urzyme (LeuAC (30uM) and GlyCA (15uM), for 30 and 60 min, respectively).
682 Unreacted amino acids were removed by running each sample through two size-exclusion columns (Zeba,
683 7K, ThermoFisher). NHS-biotin (ThermoFisher) was added in excess on ice to attach biotin to any RNA
684 with a primary amine as previously described ([Pütz, et al. 1997](#); [Chumachenko, et al. 2009](#)). Unreacted
685 NHS-biotin was removed by running each sample through two size-exclusion columns. Biotinylated RNAs
686 were captured using streptavidin magnetic beads in 0.1M sodium acetate (pH 5.2) to protect the substrate
687 from degradation. The inactive RNAs were removed by three washes of the beads with sodium acetate.
688 Biotinylated RNAs were released from the beads by breaking the RNA-amino acid bond through incubation
689 in 50mM HEPES pH 7.5 for 30 min at 37°C.

690 After the partitioning steps, short RNA oligonucleotides were ligated to the ends of the active RNAs to
691 facilitate reverse transcription and PCR needed for sequencing using the TruSeq small RNA library prep
692 kit (Illumina). We adapted the ligation steps to 24h, 15°C in the presence of 27% PEG 8000 (Promega) to
693 minimize ligation sequence bias, as previously described ([Song, et al. 2014](#)). Duplicate selections and the
694 starting library were multiplexed and sequenced on iSEQ (Illumina), yielding about 1 million sequences
695 per population. The constant regions of the sequences were trimmed off prior to analysis using our python
696 script ([Popović, et al. 2015](#)). For each population the filtered reads were counted and compared between
697 populations using FASTAptamer toolkit ([Alam, et al. 2015](#)).

698 **Identification of Class I, II bidirectional urzyme genes.**

699 The architecture of the Goldilocks LeuAC urzyme (Fig. 5D) provided the template for constructing and
700 testing possible antiparallel alignments. Pymol ([Pymol 2010](#)) representations of candidate urzymes
701 provided the platform for comparing alternative alignments between Class I and II urzymes in which we
702 deleted the second crossover connection. After transferring the sequences saved from pymol as .fasta files
703 to Excel and accessing the codon middle bases, we manually counted the number of middle base pairs in a
704 small number of relative antisense alignments. Having established a secure hypothesis from this preliminary
705 study, we assembled three alignments from aaRS.online ([Douglas, Cui, et al. 2024](#)). The amino acid
706 sequences were aligned for both candidate Classes to the secondary structures as well as the codon middle
707 bases, both of which can be represented by a single letter. We used that platform to assemble multiple
708 sequence alignments, which then were analyzed by an automated procedure to yield a matrix of middle-
709 base pairing frequencies.

710 **Data availability**

711 All data underlying this article will be shared on reasonable request to the corresponding author.

712 **Funding**

713 This work was supported by the Alfred P. Sloan Foundation Matter-to-Life Program (Grant G-2021-
714 16944); the National Institutes of Health (Grant R35GM131923 (B.K.) and NSF fellowship DGE-2040435
715 (H.D.). H.D. acknowledges support by a Pre-doctoral Fellowship from the American Foundation for
716 Pharmaceutical Education. This work utilized the resources of the UNC Longleaf high-performance
717 computing cluster.

718 **Design of bidirectional genes with Class I and II AARS urzymes on opposite strands.**

719 Initial working models for bidirectional gene construction was performed with Pymol on grounds suggested
720 by the Leucyl Goldilocks urzyme. The corresponding Class II AARS urzyme was previously thought to be

721 continuous between Motif 1 and Motif 2. The discovery that the second crossover connection of the
722 Rossmann fold is not needed for catalytic activity required some accommodations in the Class II
723 complement because of the overlap between it and the continuous sequence linking Motifs 1 and 2. These
724 accommodations were identified by anchoring the two Class I and II signatures. The Class II N-terminus
725 was then determined by the length of sequence included following the KMSKS sequence in the Class I
726 partner.

727 Structures of the TrpRS (Class I) and HisRS (Class II) enzymes were obtained, and any extant strand breaks
728 were repaired using Rosetta loop modeling (Mandell et al., 2009, <https://doi.org/10.1038/nmeth0809-551>).
729 To design bidirectional genes, we modified sequence design model ProteinMPNN (Dauparas et al.,
730 2022, <https://doi.org/10.1126/science.add2187>) to enforce bidirectional genetic coding constraints. To
731 preserve the consensus Class I and II sequence signatures, we fixed various active site residues (Class I:
732 HIGH motif, S5, Q8, P9, A21, H42, M51, D54, K81, K84; Class II: VTDV motif, R86, E87) and designed
733 all remaining positions. We generated 10,000 sequences using the default sampling temperature of 0.1.
734 Code for our modified version of ProteinMPNN capable of designing bidirectional genes can be found
735 at <https://github.com/Kuhlman-Lab/proteinmpnn>.

736 **References**

- 737
- 738 Abramson J, Adler J, Dunger J, Evans R, Green T, Pritzel A, Ronneberger O, Willmore L, Ballard AJ,
739 Bambrick J, et al. 2024. Accurate structure prediction of biomolecular interactions with AlphaFold 3. *Nature*
740 <https://doi.org/10.1038/s41586-024-07487-w>
- 741 Alam KK, Chang JL, Burke DH. 2015. FASTAptamer: A Bioinformatic Toolkit for High-throughput
742 Sequence Analysis of Combinatorial Selections. *Molecular Therapy—Nucleic Acids* 4:e230.
- 743 Ball P. 2023. *How Life Works: A User’s Guide to the New Biology*: University of Chicago.
- 744 Berg P, Ofengand EJ. 1958. An Enzymatic Mechanism for Linking Amino Acids to RNA.
745 *Proc Nat Acad Sci USA* 44:78-85.
- 746 Bokma F. 2002. Detection of punctuated equilibrium from molecular phylogenies. *Journal of Evolutionary*
747 *Biology* 15:1048–1056.
- 748 Bridgham JT, Ortlund EA, Thornton JW. 2009. An epistatic ratchet constrains the direction of
749 glucocorticoid receptor evolution. *Nature* 461:515-519.
- 750 Carter CW, Jr. 2015. What RNA World? Why a Peptide/RNA Partnership Merits Renewed Experimental
751 Attention. *Life* 5:294-320.
- 752 Carter CW, Jr, Wills PR. 2018a. Hierarchical groove discrimination by Class I and II aminoacyl-tRNA
753 synthetases reveals a palimpsest of the operational RNA code in the tRNA acceptor-stem bases. *Nucleic*
754 *Acids Research* 46:9667–9683.
- 755 Carter CW, Jr, Wills PR. 2019a. Class I and II aminoacyl-tRNA synthetase tRNA groove discrimination
756 created the first synthetase•tRNA cognate pairs and was therefore essential to the origin of genetic coding.
757 *IUBMB Life* 71:1088–1098.
- 758 Carter CW, Jr, Wills PR. 2018b. Interdependence, Reflexivity, Fidelity, and Impedance Matching, and the
759 Evolution of Genetic Coding. *Molecular Biology and Evolution* 35:269-286.
- 760 Carter CW, Jr. 2016. An Alternative to the RNA World. *Natural History* 125:28-33.
- 761 Carter CW, Jr. 2024. Base Pairing Promoted the Self-Organization of Genetic Coding, Catalysis, and Free-
762 Energy Transduction. *MDPI Life* 14:199.

- 763 Carter CW, Jr. 2022. How did the proteome emerge from pre-biotic chemistry? In: Fiore M, editor. Pre-
764 Biotic Chemistry and Life's Origin. London, UK: The Royal Society of Chemistry. p. 317-346.
- 765 Carter CW, Jr. 2014. Urzymology: Experimental Access to a Key Transition in the Appearance of Enzymes.
766 *J. Biol. Chem.* 289:30213–30220.
- 767 Carter CW, Jr. 2008. Whence the Genetic Code?: Thawing the 'Frozen Accident'. *Heredity* 100:339-340.
- 768 Carter CW, Jr., Li L, Weinreb V, Collier M, Gonzales-Rivera K, Jimenez-Rodriguez M, Erdogan O,
769 Chandrasekharan SN. 2014. The Rodin-Ohno Hypothesis That Two Enzyme Superfamilies Descended
770 from One Ancestral Gene: An Unlikely Scenario for the Origins of Translation That Will Not Be Dismissed.
771 *Biology Direct* 9:11.
- 772 Carter CW, Jr., Poppinga A, Bouckaert R, Wills PR. 2022. Multidimensional Phylogenetic Metrics Identify
773 Class I Aminoacyl-tRNA Synthetase Evolutionary Mosaicity and Inter-modular Coupling. *International*
774 *Journal of Molecular Sciences* 23: 1520.
- 775 Carter CW, Jr., Wills PR. 2019b. Experimental Solutions to Problems Defining the Origin of Codon-
776 Directed Protein Synthesis. *Biosystems* 183:103979.
- 777 Carter CW, Jr., Wills PR. 2021a. Reciprocally-coupled Gating: Strange Loops in Bioenergetics, Genetics,
778 and Catalysis. *Biomolecules* 11:265.
- 779 Carter CW, Jr., Wills PR. 2021b. The Roots of Genetic Coding in Aminoacyl-tRNA Synthetase Duality
780 *Annual Review of Biochemistry* 90:349-373.
- 781 Cestari I, Stuart K. 2013. A spectrophotometric assay for quantitative measurement of aminoacyl-tRNA
782 synthetase activity. *J Biomol Screen.* 18:490–497.
- 783 Chandrasekaran SN, Yardimci G, Erdogan O, Roach JM, Carter CW, Jr. 2013. Statistical Evaluation of the
784 Rodin-Ohno Hypothesis: Sense/Antisense Coding of Ancestral Class I and II Aminoacyl-tRNA
785 Synthetases. *Molecular Biology and Evolution* 30:1588-1604.
- 786 Chumachenko NV, Novikov Y, Yarus M. 2009. Rapid and simple ribozymic aminoacylation using three
787 conserved nucleotides. *Journal of the American Chemical Society* 131:5257-5263.
- 788 Crick FHC. 1970. Central Dogma of Molecular Biology. *Nature* 227:561-563.
- 789 Crick FHC. 1968. The Origin of the Genetic Code. *Journal of Molecular Biology* 38:367-379.
- 790 Crick FHC, Brenner, S., Klug, A., and Pieczeknik, G. 1976. A Speculation on the Origin of Protein Synthesis.
791 *Origins of Life* 7:389-397.
- 792 Crooks GE, Hon G, Chandonia J-M, Brenner SE. 2004. WebLogo: A Sequence Logo Generator. *Genome*
793 *Research* 14:1188-1190.
- 794 Cusack S. 1994. Evolutionary Implications. *Nature Structural and Molecular Biology* 1:760.
- 795 Delarue M. 2007. An asymmetric underlying rule in the assignment of codons: Possible clue to a quick
796 early evolution of the genetic code via successive binary choices. *RNA* 13:1-9.
- 797 Delarue M, Moras D. 1992. Aminoacyl-tRNA synthetases: Partition into two classes. In: Eckstein F, Lilley
798 DMJ, editors. *Nucleic Acids and Molecular Biology*. Berlin, Heidelberg: Springer-Verlag. p. 203-224.
- 799 Dill K, Chan HS. 1997. From Levinthal to pathways to funnels. *Nature Structural Biology* 4:10-19.
- 800 Dill KA, MacCallum JL. 2012. The Protein-Folding Problem, 50 Years On. *Science* 338:1042-1046.
- 801 Doublé S, Gilmore CJ, Bricogne G, Carter CW, Jr. 1995. Tryptophanyl-tRNA synthetase crystal structure
802 reveals an unexpected homology to tyrosyl-tRNA synthetase. *STRUCTURE* 3:17-31.

- 803 Douglas J, Bouckaert R, Carter CWJ, Wills P. 2024. Enzymic recognition of amino acids drove the
804 evolution of primordial genetic codes. *Nucleic Acids Research* 52:558–571.
- 805 Douglas J, Bouckaert R, Harris S, Carter CW, Jr, Wills PR. 2024. Evolution is coupled with branching
806 across many granularities of life. *Proceedings of the Royal Society B In Review*.
- 807 Douglas J, Carter CW, Jr, Wills PR. 2024. HetMM: A Michaelis-Menten model for non-homogeneous
808 enzyme mixtures. *iScience* 27:108977.
- 809 Douglas J, Cui H, Perona JJ, Vargas-Rodriguez O, Tyynismaa H, Carreño CA, Ling J, Ribas-de-Pouplana
810 L, Yang X-L, Ibba M, et al. 2024. AARS Online: a collaborative database on the structure, function, and
811 evolution of the aminoacyl-tRNA synthetases. *Life In Press*.
- 812 Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic
813 Acids Research* 32:1792-1797.
- 814 Eldredge N GSJ. 1972. Punctuated equilibria: an alternative to phyletic gradualism. In: Schopf TJM, editor.
815 *Models in paleobiology*. San Francisco: Freeman. p. 82-115.
- 816 Eriani G, Delarue M, Poch O, Gangloff J, Moras D. 1990. Partition of tRNA synthetases into two classed
817 based on mutually exclusive sets of sequence motifs. *Nature* 347:203-206.
- 818 Francklyn C, Schimmel P. 1989. Aminoacylation of RNA Minihelices with Alanine. *Nature* 337:478-481.
- 819 Fritsch EF, Sambrook J, Maniatis T. 1982. *Molecular Cloning: A Laboratory Manual*. New York.: Cold
820 Spring Harbor Laboratory.
- 821 Frugier M, Florentz C, Giegé R. 1992. Anticodon-Independent Valylation of an RNA Minihelix.
822 *Proceedings of the National Academy of Sciences USA* 89:3900-3904.
- 823 Giegé R. 1972. Recherches sur la spécificité de reconnaissance des acides ribonucléiques de transfert par
824 les aminoacyl-tRNA synthétases [Study on the specificity of recognition of transfer ribonucleic acids by
825 aminoacyl-tRNA synthetases]. [Thèse de Doctorat d'Etat]. [Strasbourg, France]: Université Louis Pasteur.
- 826 Guseva E, Zuckermann RN, Dill KA. 2017. Foldamer hypothesis for the growth and sequence
827 differentiation of prebiotic polymers. *Proc. Nat. Acad. Sci. USA* 114: E7460–E7468.
- 828 Han L, Luo Z, Ju Y, Chen B, Taotao Z, Wang J, Xu J, Gu Q, Yang X-L, Schimmel P, Zhou H. 2023. The
829 binding mode of orphan glycyl-tRNA synthetase with tRNA supports the synthetase classification and
830 reveals large domain movements. *Science Advances* 9:eadf1027.
- 831 Hoagland MB, E. B. Keller, Zamecnik. PC. 1956. Enzymatic Carboxyl Activation of Amino Acids. . *J. Biol.
832 Chem.* 21:345-358.
- 833 Hofstadter DR. 1979. Gödel, Escher, Bach: an eternal golden braid. New York: Basic Books, Inc.
- 834 Hofstadter DR. 2007. *I Am A Strange Loop*. Philadelphia, PA: Basic Books.
- 835 Hordijk W. 2019. A History of Autocatalytic Sets: A Tribute to Stuart Kauffman. *Biological Theory* 14:224–
836 246.
- 837 Jones OW, Jr. , Nirenberg MW. 1966a. Degeneracy in the amino acid code. *Biochimica et Biophysica Acta*
838 (BBA) - *Nucleic Acids and Protein Synthesis* 119:400-406.
- 839 Jones OW, Nirenberg MW. 1966b. Degeneracy in the amino acid code", . *Biochim. Biophys. Acta* 2:400–
840 406.
- 841 Katoh K, Rozewicki J, Yamada KD. 2019. MAFFT online service: multiple sequence alignment, interactive
842 sequence choice and visualization. *Brief. Bioinform.* 20:1160-1166.

- 843 Katsnelson MI, Wolf YI, Koonin EV. 2019. On the feasibility of saltational evolution. *Proc. Nat. Acad. Sci.*
844 USA 116 21068–21075.
- 845 Kauffman SA. 1986. Autocatalytic Sets of Proteins. *J. Theor. Bio.*, 119:1–24.
- 846 Klipcan L, Safto M. 2004. Amino acid biogenesis, evolution of the genetic code and aminoacyl-tRNA
847 synthetases. *Journal of Theoretical Biology* 228:389–396.
- 848 Koonin EV, Novozhilov AS. 2017. Origin and Evolution of the Universal Genetic Code. *Annu. Rev. Genet.*
849 51:45–62.
- 850 Lanier KA, Petrov AS, Williams LD. 2017. The Central Symbiosis of Molecular Biology: Molecules in
851 Mutualism. *J. Mol. Evol.* 85.
- 852 Le SQ, Gascuel O. 2008. An Improved, General Amino-Acid Replacement Matrix. *Molecular Biology and*
853 *Evolution.* 2008 25:1307-1320.
- 854 Levinthal C. 1968. ARE THERE PATHWAYS FOR PROTEIN FOLDING ? *Journal de Chimie Physique*
855 65:44.
- 856 Li L, Weinreb V, Francklyn C, Carter CW, Jr. 2011. Histidyl-tRNA Synthetase Urzymes: Class I and II
857 Aminoacyl-tRNA Synthetase Urzymes have Comparable Catalytic Activities for Cognate Amino Acid
858 Activation. *J. Biol. Chem.* 286:10387-10395.
- 859 Liu Z, Wu L-F, Xu J, Bonfio C, Russell DA, Sutherland JD. 2020. Harnessing chemical energy for the
860 activation and joining of prebiotic building blocks. *Nature Chemistry* 12:1023–1028.
- 861 Lynch M, Force A. 2000. The probability of duplicate gene preservation by subfunctionalization. *Genetics*
862 154:459–473.
- 863 Manceau M, Marin J, Morlon H, Lambert A. 2020. Model-Based Inference of Punctuated Molecular
864 Evolution. *Molecular Biology and Evolution* 37:3308–3323.
- 865 Martinez-Rodriguez L, Jimenez-Rodriguez M, Gonzalez-Rivera K, Williams T, Li L, Weinreb V,
866 Chandrasekaran SN, Collier M, Ambroggio X, Kuhlman B, et al. 2015. Functional Class I and II Amino
867 Acid Activating Enzymes Can Be Coded by Opposite Strands of the Same Gene. *J. Biol. Chem.* 290:19710–
868 19725.
- 869 Muñoz V, Serrano L. 1994. Intrinsic Secondary Structure Propensities of the Amino Acids, Using Statistical
870 f-y matrices: Comparison with Experimental Scales. *PROTEINS: Structure, Function, and Genetics*
871 20:301-311.
- 872 Newman DV. 1996. EMERGENCE AND STRANGE ATTRACTORS. *Philosophy of Science* 63:245-261.
- 873 Nirenberg MW, Matthaei JH. 1961. The Dependence of Cell-Free Protein Synthesis in *E. coli* Upon Naturally
874 Occurring or Synthetic Polyribonucleotides. *Proceedings of the National Academy of Sciences USA*
875 47:1588-1602.
- 876 Onodera K, Suganuma N, Takano H, Sugita Y, Shoji T, Minobe A, Yamaki N, Otsuka R, Mutsuro-Aoki H,
877 Umehara T, Tamura K. 2021. Amino acid activation analysis of primitive aminoacyl-tRNA synthetases
878 encoded by both strands of a single gene using the malachite green assay. *Biosystems* 208:104481.
- 879 Ostrem DL, Berg P. 1974. Glycyl Transfer Ribonucleic Acid Synthetase from *Escherichia coli*; Purification,
880 Properties, and Substrate Binding! *Biochemistry* 13:1338
- 881 Pagel M, Venditti C, Meade A. 2006. Large punctuational contribution of speciation to evolutionary
882 divergence at the molecular level. *Science* 5796:119–121.

- 883 Patra SK, Betts L, Tang GQ, Douglas J, Wills PR, Bouckear R, Carter CW, Jr. . 2024. A genomic database
884 furnishes minimal functional glycyyl-tRNA synthetases homologous to other, designed class II urzymes.
885 *Nucleic Acids Research* In Press.
- 886 Pham Y, Kuhlman B, Butterfoss GL, Hu H, Weinreb V, Carter CW, Jr. 2010. Tryptophanyl-tRNA synthetase
887 Urzyme: a model to recapitulate molecular evolution and investigate intramolecular complementation. *J.*
888 *Biol. Chem.* 285:38590-38601.
- 889 Pham Y, Li L, Kim A, Erdogan O, Weinreb V, Butterfoss G, Kuhlman B, Carter CW, Jr. 2007. A Minimal
890 TrpRS Catalytic Domain Supports Sense/Antisense Ancestry of Class I and II Aminoacyl-tRNA
891 Synthetases. *Mol. Cell* 25:851-862.
- 892 Picard V, Ersdal-Badju E, Lu A, Bock SC. 1994. A rapid and efficient one-tube PCR-based mutagenesis
893 technique using Pfu DNA polymerase. *Nucleic Acids Research* 22:2587-2591.
- 894 Popović M, Fliss PS, Ditzler MA. 2015. In vitro evolution of distinct self-cleaving ribozymes in diverse
895 environments. . *Nucleic Acids Res.* 43:7070-7082.
- 896 Pütz J, Wientges J, Schwienhorst A, Sissler M, Giegé R, Florentz C. 1997. Rapid selection of aminoacyl-
897 tRNAs based on biotinylation of α -NH₂ group of charged amino acids. *Nucleic Acids Research* 25:1862-
898 1863.
- 899 Pymol. 2010. The PyMOL Molecular Graphics System. New York, NY: Schrödinger, LLC.
- 900 Ribas de Pouplana L, Frugier M, Quinn C, Schimmel P. 1996. Evidence that two present-day components
901 needed for the genetic code appeared after nucleated cells separated from eubacteria. *Proceedings of the*
902 *National Academy of Sciences, USA* 93:166-170.
- 903 Ribas de Pouplana L, Schimmel P. 2001a. Aminoacyl-tRNA synthetases: potential markers of genetic code
904 development. *Trends in Biochemical Sciences* 26:591-596.
- 905 Ribas de Pouplana L, Schimmel P. 2001b. Two Classes of tRNA Synthetases Suggested by Sterically
906 Compatible Dockings on tRNA Acceptor Stem. *Cell* 104:191-193.
- 907 Ribas L, Schimmel P. 2001. Two Classes of tRNA Synthetases Suggested by Sterically Compatible
908 Dockings on tRNA Acceptor Stem. *Cell* 104:191-193.
- 909 Rodin SN, Ohno S. 1995. Two Types of Aminoacyl-tRNA Synthetases Could be Originally Encoded by
910 Complementary Strands of the Same Nucleic Acid. *Origins of Life and Evolution of the Biosphere* 25:565-
911 589.
- 912 Rodin SN, Rodin A. 2006. Partitioning of Aminoacyl-tRNA Synthetases in Two Classes Could Have Been
913 Encoded in a Strand-Symmetric RNA World. *DNA and Cell Biology* 25:617-626.
- 914 S. Rastogi, D. A. Liberles. 2005. Subfunctionalization of duplicated genes as a transition state to
915 neofunctionalization. *BMC Evolutionary Biology* 5: 28.
- 916 Safro M, Klipcan L. 2013. The mechanistic and evolutionary aspects of the 2'- and 3'-OH paradigm in
917 biosynthetic machinery. *Biology Direct* 8:17.
- 918 Schimmel P. 1991. RNA minihelices and the decoding of genetic information. *The FASEB Journal* 5:2180-
919 2187.
- 920 Schimmel P, Giegé R, Moras D, Yokoyama S. 1993. An operational RNA code for amino acids and possible
921 relationship to genetic code. *Proceedings of the National Academy of Sciences, USA* 90:8763-8768.
- 922 Serrano L, Sancho J, Hirshberg M, Fersht AR. 1992. Alpha-helix stability in proteins. I. Empirical
923 correlations concerning substitution of side-chains at the N and C-caps and the replacement of alanine by
924 glycine or serine at solvent-exposed surfaces. *J Mol Biol* 227:544-559.

- 925 Shore J, Holland BR, Sumner JG, Nieselt K, Wills PR. 2019. The Ancient Operational Code is Embedded
926 in the Amino Acid Substitution Matrix and aaRS Phylogenies. *Journal of Molecular Evolution* 88:136–150.
- 927 Sobotta J, Geisberger T, Moosmann C, Scheidler CM, Eisenreich W, Wächtershäuser Gn, Huber C. 2020.
928 A Possible Primordial Acetyleno/Carboxydutrophic Core Metabolism. *Life* 10:35.
- 929 Succi ND, Onuchic JN, Wolynes PG. 1998. Protein Folding Mechanisms and the Multidimensional Folding
930 Funnel. *PROTEINS: Structure, Function, and Genetics* 32:136–158.
- 931 Song Y, Liu KJ, Wang T-H. 2014. Elimination of ligation dependent artifacts in T4 RNA ligase to achieve
932 high efficiency and low bias microRNA capture. *PLoS ONE* 9:e94619.
- 933 Stubbs RT, Yadav M, Krishnamurthy R, Springsteen GR. 2020. A plausible metal-free ancestral analogue
934 of the Krebs cycle composed entirely of α -ketoacids. *Nature Chemistry* 12:1016–1022.
- 935 Takénaka A, Moras D. 2020. Correlation between equi-partition of aminoacyl-tRNA synthetases and
936 amino-acid biosynthesis pathways. *Nucleic Acids Research*, 48:3277–3285.
- 937 Takeuchi N, Hogeweg P, Kaneko K. 2017. Conceptualizing the origin of life in terms of evolution. *Phil.*
938 *Trans. R. Soc. A* 375:20160346.
- 939 Tang GQ, Elder JJH, Douglas J, Carter CW, Jr. . 2023. Domain Acquisition by Class I Aminoacyl-tRNA
940 Synthetase Urzymes Coordinated the Catalytic Functions of HVGH and KMSKS Motifs. *Nucleic Acids*
941 *Research* 51:8070–8084.
- 942 Tang GQ, Hu H, Douglas J, Carter CW, Jr. 2024. Primordial aminoacyl-tRNA synthetases preferred tRNA
943 minihelix substrates over full-length tRNA. *Nucleic Acids Research* 52:1-24.
- 944 Trupin JS, Rottman FM, Brimacombe R, Leder P, Bernfield MR, Nirenberg M. 1965. RNA Codewords and
945 Protein Synthesis, VI. On the Nucleotide Sequences of Degenerate Codeword Sets for Isoleucine, Tyrosine,
946 Asparagine, and Lysine. *Proc. Natl. Acad. Sci. U.S.A.* 53 807–811.
- 947 Watson JD, Crick FHC. 1953. A Structure for Deoxyribose Nucleic Acid. *Nature* 171:737-738.
- 948 Whelan S, Goldman N. 2001a. A General Empirical Model of Protein Evolution Derived from Multiple
949 Protein Families Using a Maximum-Likelihood Approach. *Molecular Biology and Evolution* 18:691-699.
- 950 Whelan S, Goldman N. 2001b. A General Empirical Model of Protein Evolution Derived from Multiple
951 Protein Families Using a Maximum-Likelihood Approach. *Molecular Biology and Evolution* 18:691–699.
- 952 Wills PR. 2001. Autocatalysis, information, and coding. *Biosystems* 50:49-57.
- 953 Wills PR, Carter CW, Jr. 2018. Insuperable problems of an initial genetic code emerging from an RNA
954 World. *Biosystems* 164:155-166.
- 955 Wills PR, Carter CW, Jr. 2020. Impedance matching and the choice between alternative pathways for the
956 origin of genetic coding. *International Journal of Molecular Sciences* 21:7392.
- 957 Woese CR, Olsen GJ, Ibba M, Soll D. 2000. Aminoacyl-tRNA Synthetases, the Genetic Code, and the
958 Evolutionary Process. *Microbiology and Molecular Biology Reviews* 64:202–236.
- 959 Wong JT-F, Ng S-K, Mat W-K, Hu T, Xue H. 2016. Coevolution Theory of the Genetic Code at Age Forty:
960 Pathway to Translation and Synthetic Life. *Life* 6:12.
- 961 Yarus M, Widmann J, Knight R. 2009. RNA-amino acid binding: A stereochemical era for the genetic code.
962 *Journal of Molecular Evolution* 69:406-429.
- 963 Zhu TF, Budin I, Szostak JW. 2013. Preparation of Fatty Acid or Phospholipid Vesicles by Thin-film
964 Rehydration. *Methods in Enzymology* 533:267-274.

965 Zull JE, Smith SK. 1990. Is genetic code redundancy related to retention of structural information in both
966 DNA strands? Trends in Biochemical Sciences 15:257-261.

967

968