

1 **Structured proteins are abundant in unevolved sequence space**

2

3 Vyacheslav Tretyachenko^{1,2}, Jiří Vymětal³, Tereza Neuwirthová^{1#}, Jiří Vondrášek³, Kosuke
4 Fujishima^{4,5} and Klára Hlouchová^{*1,3}

5 ¹ Department of Cell Biology, Faculty of Science, Charles University, BIOCEV, Prague, 12843, Czech
6 Republic

7 ² Department of Biochemistry, Faculty of Science, Charles University, Prague, 12843, Czech Republic

8 ³ Institute of Organic Chemistry and Biochemistry, Czech Academy of Sciences, Prague, 16610, Czech
9 Republic

10 ⁴ Earth-Life Science Institute, Tokyo Institute of Technology, Tokyo, 1528550, Japan

11 ⁵ Graduate School of Media and Governance, Keio University, Fujisawa, Japan, 2520882

12 # Current address: R&D Informatics Solutions, MSD Czech Republic s.r.o., Prague, Czech Republic

13 * To whom correspondence may be addressed: klara.hlouchova@natur.cuni.cz

14 **Abstract**

15 Natural proteins represent numerous but tiny structure/function islands in a vast ocean of possible
16 protein sequences, most of which has not been explored by either biological evolution or research.
17 Recent studies have suggested this uncharted sequence space possesses surprisingly high
18 structural propensity, but development of an understanding of this phenomenon has been
19 awaiting a systematic high-throughput approach.

20 Here, we designed, prepared, and characterized two combinatorial protein libraries consisting of
21 randomized proteins, each 105 residues in length. The first library constructed proteins from the
22 entire canonical alphabet of 20 amino acids. The second library used a subset of only 10 residues
23 (A,S,D,G,L,I,P,T,E,V) that represent a consensus view of plausibly available amino acids through
24 prebiotic chemistry. Our study shows that compact structure occurrence (i) is abundant (up to
25 40%) in random sequence space, (ii) is independent of general Hsp70 chaperone system activity,
26 and (iii) is not granted solely by “late” and complex amino acid additions. The Hsp70 chaperone
27 system effectively increases solubility and stability of the canonical alphabet but has only a minor
28 impact on the “early” library. The early alphabet proteins are inherently more stable and soluble,
29 possibly assisted by salts and cofactors in the cell-like environment in which these assays were
30 performed.

31 Our work indicates that natural protein space may have been selected to some extent by chance
32 rather than unique structural characteristics.

33

34 **Keywords**

35 Protein sequence space, protein structure, amino acid alphabet, genetic code evolution, random
36 proteins

1 Introduction

2 Today's biological systems are anchored in the universal genetic coding apparatus, relying on
3 coded amino acids that were likely selected in the first 10-15% of Earth's history ¹. While sources
4 of prebiotic organic material provided a broad selection of amino acids, only about half of the
5 canonical amino acids were detected in this pool ². There is substantial evidence that this set
6 formed an early version of the genetic code and that the "late" amino acids were recruited only
7 after an early metabolism was in existence. The boundary between these two sets is blurry.
8 However, large meta-analyses of these studies agree that "early", i.e. the smaller and less
9 complex amino acids (Gly, Ala, Asp, Glu, Val, Ser, Ile, Leu, Pro, Thr) were a fixture in the genetic
10 code before its evolution ^{3,4}.

11 The factors that drove the selection of 20 coded amino acids remain puzzling. Solubility, ease of
12 biosynthesis, un/reactivity with tRNA, and potential peptide product stability seem to explain some
13 selective "choices" but not others ^{5,6}. Most recently, analysis of the *set* of amino acids revealed
14 that the canonical alphabet shows an unusually good repertoire of the chemical property space
15 when compared to plausible alternatives ^{7,8}. Such studies lead to speculations that similar amino
16 acid selection would be expected on other Earth-like planets ^{5,8,9}.

17 In extant proteins, a significant part of the "late" amino acids (Arg, Lys, His, Cys, Trp and Tyr)
18 belong to the essential catalytic residues, i.e. they are associated with catalysis in almost all of
19 the enzyme classes ¹⁰. At the same time, the putatively early amino acids have been related to
20 protein disorder and lack of 3D structure ¹¹. However, scarce sampling of random sequences
21 composed of early amino acids suggests that such proteins have a higher solubility than the full
22 canonical alphabet ^{12,13}. Moreover, computational and experimental mutational studies removing
23 or reducing the late amino acids in selected proteins imply that the early amino acids comprise a
24 non-zero folding potential ¹⁴⁻¹⁸. If prone to tertiary structure formation, it has been hypothesized
25 that the early alphabet could more probably form molten globules rather than tightly packed
26 structures, mainly due to the lack of aromatic and positively charged amino acids. According to
27 this hypothesis, the addition of late amino acids would be required to increase protein stability
28 and catalytic activity ^{11,17,19}. Interestingly, it was shown that while positively charged amino acids
29 are more compatible with protein folding, they also promote protein aggregation if their position
30 within the sequence is not optimized or assisted by molecular chaperones. Thus it was
31 hypothesized that chaperone emergence coincided with the incorporation of basic residues into
32 the amino acid alphabet leading to an increase in the plasticity of natural folding space ²⁰.

33 To assess the intrinsic structural and functional properties of the full amino acid alphabet, semi
34 high-throughput studies using combinatorial sequence libraries have been performed previously
35 ²¹⁻²⁵. Surprisingly, secondary structure occurrence in random sequence libraries has been
36 recorded with similar frequency as in biological proteins, while folding (or more precisely,
37 occurrence of collapsed conformations) has been reported in up to 20% of tested proteins ^{21,24,25}.
38 However, more systematic and high-throughput screening is still necessary to confirm these
39 observations. Moreover, it remains unclear how much these properties are a result of the full
40 alphabet fine-tuning, whether structured molecules emerge spontaneously and independently in
41 the canonical amino acid sequence space, and whether the early amino acids could provide
42 similar structural traits.

43 To fill this knowledge gap, we characterized libraries of 10^{12} randomized protein sequences from
44 the full and early amino acid alphabets to assess their collective biochemical characteristics.
45 While the bioinformatic prediction revealed similar secondary structure potential in both libraries
46 and lower aggregation propensity of the full alphabet, the early alphabet is significantly more

1 soluble and thermostable under cell-like experimental conditions. The full alphabet sequences
2 were found to interact with molecular chaperones that can compensate for their otherwise poor
3 solubility. Up to ~40% folding occurrence is observed in both studied libraries. The results
4 therefore agree with previous scarce sampling observations, and in addition, the folding frequency
5 and inducibility of some properties in a cell-like environment are systematically mapped.
6 Moreover, this study provides a unique synthetic biology pipeline that could be used to survey
7 properties of any other protein alphabets associated with different biological phenomena of
8 interest.

9

10 Results

11 *Library expression and quality control*

12 The combinatorial protein libraries studied in this work consisted of 105 amino acid long proteins
13 with an 84 amino acid long variable parts, FLAG/HIS tag sequences on N'/C' ends, and a thrombin
14 cleavage site in the middle of the protein construct (Supplementary Fig. S1). The variable region
15 was designed by the CoLiDe algorithm and consisted of a specific set of degenerate codons in
16 order to match the natural canonical (full alphabet, 20F) and the prebiotically plausible
17 (A,S,D,G,L,I,P,T,E,V; early alphabet, 10E) amino acid distributions (Supplementary Table S1)²⁶.
18 The amino acid ratios for both libraries corresponded to natural amino acid distribution from the
19 UniProt database²⁷. The libraries were assembled from two overlapping oligonucleotides,
20 transcribed into their corresponding mRNA, and translated using an in vitro translation system
21 (Supplementary Fig. S2). In order to verify the designed library variability and amino acid
22 distribution, we sequenced the assembled degenerate oligonucleotide DNA library and performed
23 a mass spectrometric analysis of the purified library protein product. The root mean squared error
24 (RMSE) from the target amino acid distribution was ~0.06 in both libraries 20F and 10E
25 (Supplementary Table S2, Supplementary Fig. S3). The variability analysis of the sequenced
26 library showed that 96% of sequences were unique; no significant sequence enrichment was
27 observed (Fig. 1, Supplementary Table S3). Due to synthesis errors, STOP codons were
28 introduced into 12% of the library sequences. However, their products were not observed in
29 western blot protein analyses (Supplementary Fig. 5/7/9). The variability of the purified protein
30 product was validated by MALDI-TOF mass spectrometry; the mean and spread of the
31 experimental spectra closely matching the predicted distributions (Supplementary Fig. S4).

32



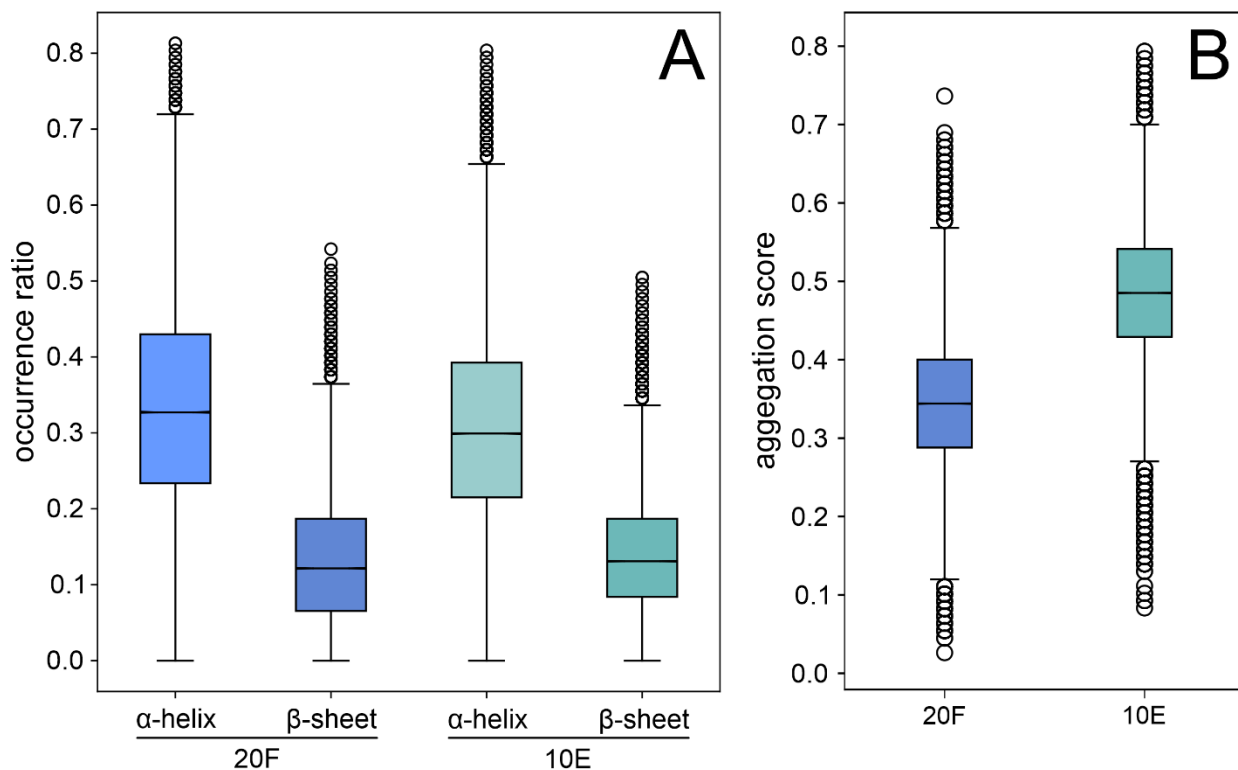
33
34 *Figure 1. Sequence logo representation of full (top) and early (bottom) alphabet libraries variability*
35 *constructed from the corresponding sequenced DNA templates.*

36

1 **Secondary structure and aggregation propensity predictions**

2 Sequences of both 20F and 10E libraries acquired by high throughput sequencing were analyzed
3 by a consensus protein secondary structure prediction²⁸. 200,000 sequences were analyzed
4 from each library. Interestingly, despite the different amino acid distributions, comparable alpha
5 helix and beta sheet forming tendencies were reported in both libraries with only a slight increase
6 in alpha helix content in the 20F library (33 % vs. 30% in 10E) (Fig. 2A). The overall alpha helix
7 and beta sheet content correlate well among the individual predictors used for both studied
8 libraries, which is not necessarily the case for other alternative and more artificial alphabets
9 (unpublished observation). The prediction of aggregation propensity of the same set of sequences
10 indicated significantly higher aggregation tendency of 10E library proteins in comparison to 20F
11 library proteins (Fig. 2B).

12



13

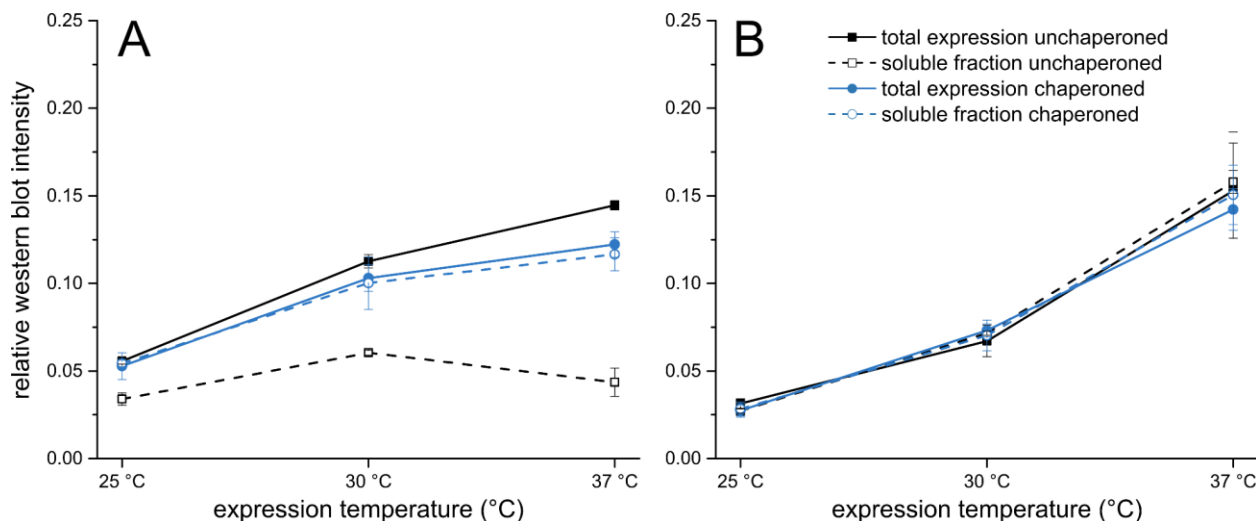
14 *Figure 2. Bioinformatic prediction of alpha helix and beta sheet content (A) and aggregation*
15 *propensity (B) of a sample of 200,000 sequences acquired by high throughput sequencing of*
16 *the early (green) and full (blue) alphabet library DNA templates. Aggregation score is defined as*
17 *the ratio of predicted aggregation-prone residues per sequence*

18 **Expression and solubility analysis in the absence and presence of the DnaK chaperone** 19 **system**

20 To systematically assess the expression profiles of the libraries, a quantitative western blot
21 analysis was performed with the library products expressed at different temperatures (25, 30 and
22 37 °C) and with/without DnaK/DnaJ/GrpE chaperone system supplementation (further referred as
23 to DnaK). The analysis was carried out in triplicate, and western blot signals of both total
24 expression and soluble fractions were quantified with ImageJ²⁹. For both 20F and 10E libraries,

1 the expression yields grow with increasing temperature, with the overall yield being mildly lower
2 in the chaperone supplemented reactions at 37 °C (Fig. 3). In the case of the 20F library, the
3 solubility of the library is relatively poor but is significantly improved by chaperone
4 supplementation. While in the chaperone supplemented reaction the soluble fraction grew with
5 expression temperature proportionally with the total expression, in the chaperone absent
6 condition, the soluble fraction yields did not significantly change with the transition from 30 to 37
7 °C. On the other hand, chaperone supplementation did not have a significant effect on the 10E
8 library expression or solubility (Fig. 3).

9



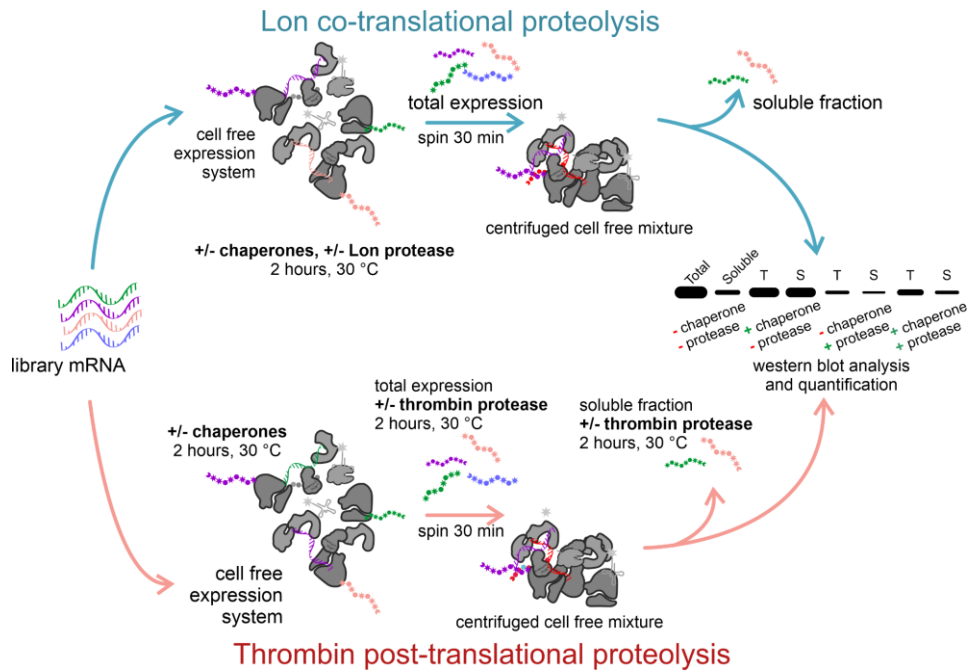
10

11 *Figure 3. A summary of expression and solubility analysis of the full (A) and early (B) alphabet*
12 *libraries at three different temperatures. Total expression (solid line) and soluble fraction (dashed*
13 *line) were compared in chaperoned (blue line) and unchaperoned (black line) conditions. For*
14 *original data see Supplementary Fig. S5/S6 and Supplementary Table S4.*

15

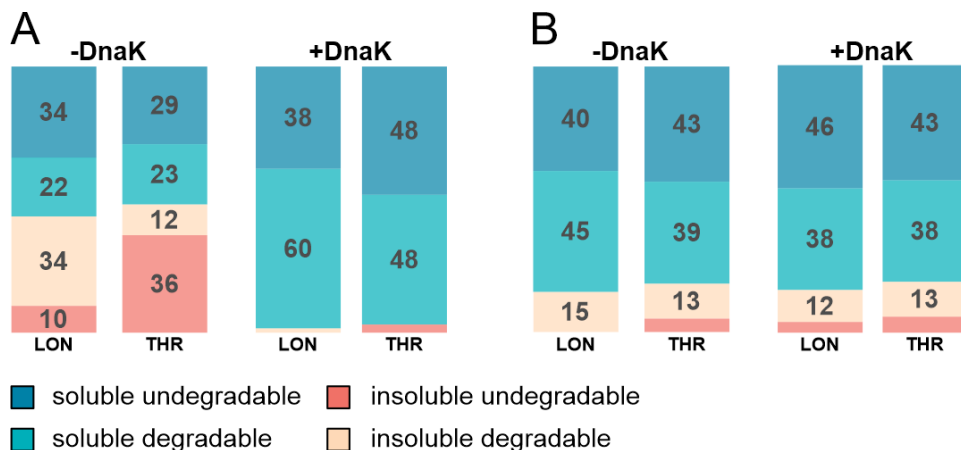
16 **Assessment of proteolytic resistance**

17 The structural potential of random protein libraries was assessed by proteolysis. The digestion
18 assessment was performed in triplicate by Lon and thrombin proteases in co-translational and
19 post-translational conditions, respectively (Fig. 4). The Lon protease is a part of the *E. coli* protein
20 misfolding system and is known to specifically digest unfolded proteins in exposed hydrophobic
21 regions³⁰. Here we adapted a previously published protocol on single protein structure
22 assessment for combinatorial library characterization³¹. The method is used to separate and
23 quantify distinct protease sensitive parts of the library within both the soluble and insoluble
24 fractions of the expressed libraries. The thrombin protease assay was adapted from the study of
25 Chiarabelli et al, wherein the structure occurrence is derived from the cleaved/uncleaved ratio of
26 proteins with an engineered thrombin cleavage site situated in the middle of the sequence²¹. The
27 unstructured proteins are expected to be quickly degraded on the exposed cleavage site.



1
2 *Figure 4. Scheme of the proteolytic resistance experimental pipeline. In the co-translational*
3 *proteolytic assay (top) the Lon protease is present during the cell-free expression; in the post-*
4 *translational proteolytic assay (bottom) thrombin protease is added to the separated total and*
5 *soluble fractions of the expressed library after translation is quenched by addition of puromycin*

6 According to the 20F library analysis, the soluble/undegradable structured proteins represent ~30-
7 35% of the total product (Fig. 5A). Upon addition of the DnaK chaperone, most of the library
8 solubilizes, but the structured content does not increase significantly and occupies ~40-50% of
9 the total product. In comparison, chaperone addition does not have an impact on the solubility or
10 structure content in the 10E library (Fig. 5B). Interestingly, the structured content (soluble
11 undegradable) in the 10E library is similar as in the 20F library after the addition of chaperones.



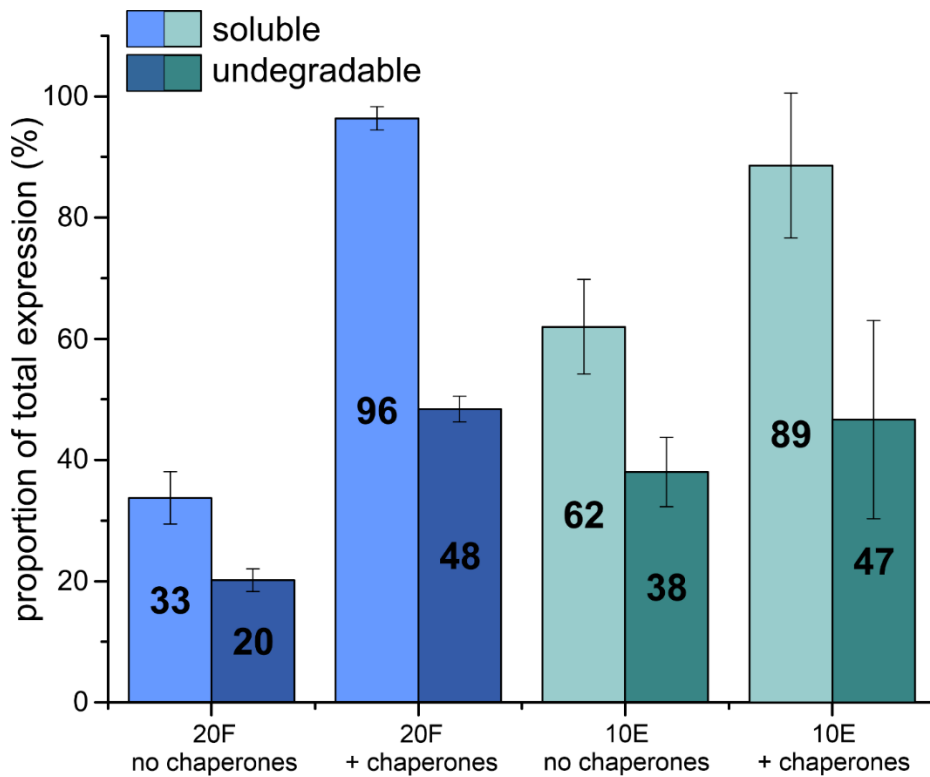
12
13 *Figure 5. An integrated solubility/proteolysis resistance analysis of the full (A) and early (B)*
14 *alphabet libraries. Libraries were expressed either in the absence (left double column) or*
15 *presence (right double column) of the DnaK chaperone system. Proteolysis was performed by*
16 *protease Lon (left columns) in a co-translational regime or by thrombin protease (right columns)*

1 *in a post-translational mode. Values in the boxes represent the percentage ratios of the total*
2 *expressed library per fraction. For original data see Supplementary Fig. S7/S8/S9/S10 and*
3 *Supplementary Table S5/S6.*

4 **Thermostability characterization**

5 Following expression, solubility, and structural content assessment, we analyzed the temperature
6 sensitivity of the 20F and 10E proteins. The libraries expressed with and without chaperone
7 supplementation were subjected to 15 minutes/42 °C heat shock. The aggregated fraction was
8 removed by centrifugation, and the soluble fraction was compared with and without thrombin
9 treatment (Fig. 6).

10



11

12 *Figure 6. Thermostability analysis showing soluble proportions (light blue and green) of the total*
13 *expression of the full and early alphabet libraries after a heat shock (42 °C / 15 min) treatment*
14 *and their respective thrombin resistant proportions (dark blue and green) of the total expression*
15 *in unchaperoned and chaperoned conditions. Numbers in the bars represent the percentage*
16 *fraction of the total expressed library. For original data see Supplementary Fig. S9/S10 and*
17 *Supplementary Table S6.*

18 The 10E library is intrinsically more thermostable than 20F (~60 vs ~30% of the libraries remain
19 soluble after heat shock, respectively) while the DnaK chaperone system induces thermostability
20 in both. The protease resistant fraction of the soluble part of the libraries remains the same (~40%)
21 as before heat shock treatment with the exception of the unchaperoned 20F library, which
22 demonstrates a slight decrease in both the soluble and degradation resistant fractions (Fig. 6).

1 Discussion

2
3 In this study, a high-throughput systematic approach was used to experimentally analyze the
4 structural properties of the vast protein sequence space. Random sequences have been
5 proposed as proxies for both (i) precursors of *de novo* emerged proteins in current evolution as
6 well as (ii) sources of peptide/protein birth at the earliest stages of life preceding templated
7 proteosynthesis^{32,33}. However, the structural properties of random sequences have so far
8 remained uncomprehended, while a few recent bioinformatic and coarse-grained studies have
9 pointed to their surprising properties, such as high secondary structure propensity and *in vivo*
10 tolerance^{24,25,34}. Here, two combinatorial protein libraries encompassing upto 10¹² individual
11 sequences from two distinct alphabets (representing hypothetical stages of genetic code
12 evolution) have been characterized.

14 ***Solubility of the natural alphabet random proteins can be induced by chaperones***

15 The first “full” alphabet library is based on the amino acid composition of the Uniprot database
16 representing the properties of today’s proteomes. It has previously been shown that similar
17 constructs have limited solubility but a similar secondary structure potential to biological proteins
18^{12,13,25}. Our study confirms these results, and in addition, we specify that 20-50% of the overall
19 diverse library appears in the soluble fraction in the 30-37 °C temperature range. While previous
20 studies of similar construct size evaluated the solubility of individual proteins that were
21 overexpressed (many of them with partial solubility) in different *E. coli* strains and under different
22 conditions, our library was expressed using a reconstituted cell-free protein synthesis (CFPS)
23 system, and its large diversity (contrasting with overexpression of individual proteins) was
24 confirmed by MALDI. Therefore, we cannot make a direct comparison to previous studies of
25 individual proteins but rather report the “fingerprint” properties of the full alphabet domain-size
26 proteins.

27 Interestingly, this library of unevolved sequences was observed to interact productively with the
28 natural molecular chaperone system DnaK/DnaJ/GrpE which was used to supplement the CFPS
29 system in another experiment. This interaction caused almost total solubilization of the otherwise
30 insoluble proteins over the studied temperature range. While the solubility traits may be quite
31 different for much shorter polymer lengths, our previous study showed that random domain-size
32 sequences cope with significant aggregation, especially if they are rich in secondary structure
33 content²⁵. To characterize the library folding potential without introducing potential bias, we used
34 an *in situ* double proteolysis experiment adapting two previously reported approaches^{21,31}. The
35 experiment combined co-translational proteolysis by disorder-specific Lon protease and a post-
36 translational cleavage by thrombin designed to cut the potentially exposed cleavage site
37 engineered in the center of random proteins. Besides the increased robustness of the structure
38 content estimation, such a combined approach provides unique insight into the library translation
39 dynamics.

40 The double proteolysis experiment revealed that ~30-35% of library 20F proteins are protease
41 resistant during proteolysis. Upon the addition of chaperones (which solubilizes the library as
42 described above), the ratio of protease resistant species rose only mildly to ~40-50%. The
43 preferentially unstructured nature of the full alphabet library echoes the nature of naturally evolved
44 *de novo* proteins, i.e. proteins that emerge in current biology from previously non-coding DNA
45 (summarized in³⁵).

46 Overall, these results show that while inherent protein solubility is limited in random sequence
47 space made of the natural alphabet, it can be induced significantly by the activity of molecular
48 chaperones. At the same time, the DnaK chaperone system has only a minor effect on structure
49 formation, suggesting that the majority of the potentially solubilized sequences are devoid of
50 tertiary structure arrangements. Nevertheless, the ~40% natural abundance of soluble and yet
51 protease-resistant sequences in unevolved sequence space may be surprising in light of earlier

1 hypotheses and exceeds the estimates of folding frequency reported by previous coarse-grained
2 studies^{21,36}. Nevertheless, major differences in the experimental setups (cell-free vs cell-based
3 expression, low-level vs overexpression, high- vs. low-throughput methodology, library design
4 and sequence length) prevent the possibility of direct comparisons among these studies. A direct
5 comparison of the full library properties can however be made with another library of proteins
6 studied here under the same experimental conditions.

7
8
9 ***Structure formation is comparable in proteins from the full canonical alphabet and its early
10 subset, unaffected by chaperones***

11 A second “early” alphabet library was constructed from a 10 amino acid subset of the full alphabet
12 which was proposed to constitute an earlier version of the genetic code and be reflected in the
13 composition of early proteins³. We emphasize that with this study, we do not try to establish that
14 there was necessarily a time in life’s evolution during which domain-size proteins were composed
15 entirely of this amino acid subset. Our analysis rather deals with the inherent physico-chemical
16 properties of such an alphabet, were it to form or dominate protein-like structures. We also
17 acknowledge that the earliest stages of peptide/protein formation (preceding templated
18 proteosynthesis and perhaps also its early less specific versions) probably utilized a plethora of
19 prebiotically plausible amino acids or similar chemical entities, but inclusion of such non-canonical
20 amino acids in the studied alphabets is beyond the scope of this study^{1,37,38}.

21 Although the overall secondary structure propensity of the early alphabet is comparable to the full
22 alphabet, according to the bioinformatic prediction, the occurrence of alpha-helix is slightly (~3%)
23 lower. While these differences are statistically borderline, they may have interesting implications
24 for the evolution of protein structural properties. Brack and Orgel proposed that beta-sheet
25 structures were prebiotically significant, and the later significance of alpha-helices in protein folds
26 was also recently implied by the structural analysis of ribosomal protein content, showing that the
27 most ancient protein-protein fragments of this molecular fossil are mostly disordered and of beta-
28 sheet formation³⁹⁻⁴¹. Despite the similar secondary structure propensities of the full and early
29 alphabets, the 10E library proteins are significantly more soluble (~90%) upon expression. They
30 retain similar solubilities in chaperoned/unchaperoned conditions unlike the 20F library proteins.
31 This observation supports the previously stated hypothesis of chaperone co-evolution with the
32 incorporation of the first positively charged amino acids into the early amino acid alphabets²⁰.

33 The significantly higher solubility of the 10E library proteins (and similar protein compositions) is
34 in agreement with previous studies^{12,13}. This phenomenon could be related to the lower
35 complexity of 10E library proteins resulting the limited amino acid alphabet. While 20F proteins
36 represent a highly variable sample of protein folding space with many opportunities for
37 aggregation initiation, the 10E proteins display a narrower subspace with much more uniform
38 sequence and physicochemical characteristic distributions. In addition, their overall negative
39 charge and absence of positively charged/aromatic amino acids are conditions which were
40 previously shown to suppress both nonspecific aggregations as well as independent protein
41 folding formation²⁰. At the same time though, the 10E alphabet contains a significant proportion
42 of hydrophobic amino acids. Using the ProA bioinformatic predictor of protein aggregations, the
43 10E library would be expected to be intrinsically less soluble, contradicting our observations as
44 well as previous empirical observations. However, contrasting with the intrinsic behavior of the
45 protein alone, our assays (and previous experimental assays) were performed in a cell-like
46 environment, rich in different salts and other small molecules/cofactors.

47 Interestingly, the 10E library also displays a significant amount of tertiary structure
48 formation. In the absence of chaperones, the ratio of the protease resistant fraction is 40-50% in
49 both the co- and post-translational digestion assay, i.e. similar to the 20F protease resistant
50 fraction when supplemented with chaperones.

1 Such a high occurrence of structure formation within the 10E library is non-intuitive and
2 unexpected purely from its amino acid composition. However, several folders have been recently
3 identified from the same or similar protein composition in experiments reducing extant protein
4 compositions^{15,16,18,42,43}. Where characterized in more detail, assistance of salts, metal ions, or
5 cofactor binding were found to explain the folding properties^{15,18,42,44}. In addition, Despotovic et
6 al. recently confirmed that folded conformations of a highly acidic 60-residue protein can be
7 induced by positively charged counterions, in case of Mg²⁺ the reported concentration
8 corresponding roughly to its concentration in the CFPS reaction (~10mM)⁴⁵. These studies allow
9 us to hypothesize that the high structural propensity of the 10E alphabet could result from the
10 cation/cofactor-rich environment, where the lack of hydrophobic and electrostatic interactions is
11 compensated by these chemical entities. Alternatively or concurrently, the library solubility and
12 protease resistance could be partly explained by tertiary structure formation induced by
13 oligomerization as previously hypothesized by Yadid et al. in a study using 100 amino acid long
14 fragments (albeit from different amino acid compositions)⁴⁶. Our study presented here cannot
15 unambiguously differentiate between these two possible scenarios or their combination as the
16 highly variable library sample of a limited amount prevents more sophisticated physico-chemical
17 analyses that could be used to address these phenomena in follow-up studies.

18

19 ***Early alphabet proteins are inherently more thermostable in a cell-like milieu***

20 One of the notable assumed characteristics of the early prebiotic Earth is the elevated
21 temperature of the environment⁴⁷. The temperature-induced aggregation propensity of random
22 protein libraries was investigated by their exposure to a 15-min heat shock at 42 °C. Interestingly,
23 the quantity of thermostable fractions in proteins without chaperones were approximately two
24 times great in the early alphabet library (~30% vs ~60% for 20F and 10E libraries, respectively)
25 which might indicate a natural tendency to withstand elevated temperature. On the other hand,
26 addition of chaperones improved the thermal stabilities of both 20F and 10E libraries up to almost
27 full solubility upon heat shock treatment. This observation confirms our previous conclusions
28 about the strong dependence of the canonical amino acid alphabet proteins on chaperone activity
29 and extends it to stability support of the early amino acid alphabet proteins. Additionally, the
30 fraction of protease resistant proteins remains unchanged (~40%) upon heat shock for both
31 libraries, suggesting that the proteins destabilized by elevated temperature do not belong to this
32 category.

33 While most of the above referenced studies reducing the composition of extant proteins towards
34 the early set of amino acids did not observe an increase in their thermostability^{15,16,18,42,44}, we are
35 here concerned with a comparison of unevolved sequences from the two amino acid repertoires
36 and their inherent properties. Unlike the studies performed with purified protein samples, our
37 thermostability assay was performed in the CFPS reaction milieu, i.e. in an environment rich in
38 salts and other small molecules, indicating innate thermostability properties in the presence of
39 such chemical entities.

40

41 ***Concluding remarks***

42 In summary, while our study confirms some of the previously reported properties of the
43 random sequences space (such as its surprisingly high secondary structure potential and relative
44 ease of expression), we expand on this knowledge using a systematic high-throughput approach
45 using diverse combinatorial libraries composed of two different alphabets. Escaping the restraints
46 of scarce sampling, our study maps the tertiary structure, solubility, and thermostability potential
47 in random sequences composed of the natural vs. the early evolutionary canonical alphabets.
48 The analyses were performed in a cell-like environment (rich in salts and cofactors) that may
49 better represent protein formation conditions during both the origins of life and in extant biology.
50 Under such conditions, the early alphabet sequences are inherently more soluble and
51 thermostable while the natural alphabet proteins can reach similar properties through interactions

1 with natural chaperones. Interestingly, our study reports that both alphabets frequently give rise
2 to proteolysis resistant soluble structures, occupying up to ~40% of all sequences. Because the
3 intrinsic properties of the prebiotically plausible amino acids do not imply such properties, we
4 hypothesize that the protein solubility and folding within this library are enabled by the cell-like
5 milieu, assisted by salts, metal cations, and cofactors. Follow up studies are suggested to further
6 explore these findings.

7 **Methods**

8 ***Design of libraries from early and full amino acid alphabet***

9 Two 105 amino acid long random sequence libraries were designed using the CoLiDe algorithm
10 for combinatorial library design²⁶ and the amino acid ratios listed in Supplementary Table S1.
11 The randomized part of the libraries consisted of 84 amino acids; the remainder is attributed to
12 the FLAG affinity purification site on the N-end of the construct, the hexahistidine tag on the C-
13 end, and the and thrombin protease recognition site (ALVPRGS) in the middle of the construct
14 (Supplementary Figure S1).

15 ***Bioinformatic analysis of secondary structure potential***

16 Prediction of secondary structure potential of the studied libraries was performed by a consensus
17 predictor as described previously²⁸. It combines outputs of the spider3, psipred, predator, jnet,
18 simpa, and GOR IV secondary structure predictors⁴⁸⁻⁵³. None of the predictors were allowed to
19 use homology information that might prevent high-throughput processing of protein sequences.
20 The final assignment of secondary structure followed the most frequently predicted secondary
21 structure element at each amino acid position. Protein aggregation was predicted by the ProA
22 algorithm in a protein prediction mode⁵⁴.

23 ***Preparation of experimental libraries***

24 20F and 10E DNA libraries were synthesized commercially as two overlapping degenerate
25 oligonucleotides (see Supplementary information for the sequences) that were designed by the
26 CoLiDe algorithm to follow the natural canonical (full alphabet, 20F) and prebiotically plausible
27 (A,S,D,G,L,I,P,T,E,V; early alphabet, 10E) amino acid distributions (Supplementary Table S1).
28 The overlapping oligonucleotides were annealed and extended by Klenow fragment to form
29 double-stranded DNA (dsDNA). Annealing was performed by heating the complementary
30 oligonucleotide mixture (48 µl total reaction volume, 2 µM final concentration of each) in NEB2
31 buffer provided with 200 µM dNTPs to 90 °C for 2 minutes and cooling down to 32 °C with a 1
32 °C/min temperature gradient. The Klenow extension was performed by Klenow polymerase
33 (NEB): 10 U of Klenow polymerase was added to annealed oligonucleotides, incubated for 5
34 minutes at 25 °C, 37°C for 1 hour (polymerization step), and 50 °C for 15 minutes (inactivation
35 step). Final dsDNA libraries were further column purified using the DNA Clean and Concentrator
36 kit (Zymo Research), and the product was quantified by Nanodrop 2000c (Thermo Scientific). In
37 the following transcription, 1 µg of DNA library was used as a template for mRNA synthesis by
38 HiScribe T7 kit (NEB). The product was purified by NH₄Ac precipitation and dissolved in RNase-
39 free water to a final concentration of 3 µg/ul.

40 The library DNA was analyzed by high throughput sequencing on Illumina MiSeq. The libraries
41 for next generation sequencing (NGS) were prepared from 100 ng DNA samples using the
42 NEBNext Ultra II DNA Library Prep kit (New England Biolabs) with AMPure XP purification beads
43 (Beckman Coulter). the length of the prepared library was determined by Agilent 2100 Bioanalyzer

1 (Agilent Technologies) and quantified by Quantus Fluorometer (Promega). The sample was
2 sequenced on a MiSeq Illumina platform using the Miseq Reagent Kit v2 500-cycles (2x250) in a
3 paired-end mode. Raw data was processed with the Galaxy platform, and sequence analysis of
4 assembled and filtered paired reads was performed with MatLab scripts developed at Heinis lab
5 ^{55,56}.

6 the protein library was expressed using the PUREflex 2.0 (GeneFrontier Corporation)
7 recombinant in vitro translation system. The reaction was supplemented by 0.05 % (v/v) Triton X-
8 100 and prepared according to manufacturer recommendations. The reaction was initiated by 3
9 µg of library mRNA. Expression followed for 2 hours at 25, 30, or 37 °C.

10

11 ***Affinity purification of protein libraries***

12

13 Expressed protein libraries were diluted 10x in binding buffer (50mM Tris, 150 mM NaCl, 0.05%
14 (v/v) Triton X-100, pH 7.5) and incubated for 2 hours at 25 °C with 3 µl / 20 µl reaction of
15 TALON affinity purification matrix. The immobilized library was washed three times with binding
16 buffer and eluted by addition of 20 µl / 20 µl reaction of elution buffer (50mM Tris, 150 mM NaCl,
17 10mM EDTA, 0.05% (v/v) Triton X-100, pH 7.5).

18

19

20 ***Solubility analysis of protein libraries***

21

22 Cell free protein expression reactions were supplemented with 0.05 % Triton X-100, and protein
23 libraries were expressed in different temperatures according to manufacturer recommendations.
24 In order to analyze the quantity of total protein product, 10 µl of each reaction was quenched by
25 addition of 40 µl of 300 µM puromycin in 50 mM Tris, 100 mM NaCl, 100 mM KCl, pH 7.5.
26 Quenching proceeded for 30 minutes at 30 °C. Next, 5 µl of the quenched reaction mixture was
27 taken for the following SDS-PAGE analysis of total library expression; the rest of the mixture was
28 centrifuged for 30 minutes at 21 °C, and 5 µl of supernatant was taken for SDS-PAGE analysis
29 of the soluble fraction of the library. Both fractions were analyzed by quantitative Western blotting
30 (Sigma-Aldrich Monoclonal ANTI-FLAG® M2-Peroxidase (HRP) antibody, A8592) following the
31 SDS-PAGE separation.

32 ***Lon proteolytic assay of protein libraries***

33

34 Lon protease was expressed and purified according to the previously published protocol ³¹. Cell
35 free expression reactions were supplemented with 0.05 % Triton X-100; reactions were prepared
36 according to manufacturer recommendations. Libraries were expressed in the presence or
37 absence of the DnaK chaperone (K+/K-) and in the presence or absence of Lon protease (L+/L-
38). Chaperones were added to the final concentration of 5 µM DnaK, 1 µM DnaJ, 1 µM GrpE and
39 Lon protease to 0.4 µM (hexamer)/reaction. Expression proceeded in 10 µl reaction volume for 2
40 hours at 30 °C and was quenched by 40 µl addition of 300 µM puromycin in 50 mM Tris, 100 mM
41 NaCl, 100 mM KCl, pH 7.5. Quenching proceeded for 30 minutes at 30 °C. The sample
42 preparation of total and soluble library fractions was identical to the solubility analysis experiment
43 described above.

44

45 ***Thrombin proteolytic assay of protein libraries***

46

47 Cell free expression reactions were supplemented with 0.05 % Triton X-100; reactions were
48 prepared according to manufacturer recommendations. Libraries were expressed in the presence

1 or absence of the chaperone DnaK (K+/K-). Chaperones were added to the final concentration of
2 5 μ M DnaK, 1 μ M DnaJ, 1 μ M GrpE μ M. Expression proceeded in 10 μ l reaction volume for 2
3 hours at 30 °C and was quenched by 40 μ l addition of 300 μ M puromycin in 50 mM Tris, 100 mM
4 NaCl, 100 mM KCl, pH 7.5. Quenching proceeded for 30 minutes at 30 °C. Post-translational
5 thrombin proteolysis was prepared as follows: 5 μ l of quenched reaction was diluted 4x by 15 μ l
6 of 50 mM Tris, 100 mM NaCl, 100 mM KCl, pH 7.5; 0.15 U of thrombin (SigmaAldrich, USA) was
7 added, and the total expressed library was digested for 2 hours at 30 °C. The soluble fraction of
8 the library was prepared by centrifugation at 21 000 xg for 30 minutes at 21 °C, and 5 μ l of
9 supernatant was thrombin digested according to the same protocol. Cleaved samples of the total
10 expressed and soluble libraries were analyzed by SDS-PAGE and Western blotting (Sigma-
11 Aldrich Monoclonal ANTI-FLAG® M2-Peroxidase (HRP) antibody, A8592).

12 ***Thermostability assay***

13 Libraries expressed in 10 μ l volume were processed as described above in the Lon proteolytic
14 assay protocol. The Lon absent libraries were further analyzed for their thermostability in the
15 presence and absence of chaperone. Processed reactions were incubated at 42 °C for 15
16 minutes and immediately centrifuged at 21 000 xg for 30 minutes at 21 °C. The 5 μ l supernatant
17 fractions were subjected to thrombin proteolysis as described previously and analyzed by SDS-
18 PAGE and quantitative western blotting.

19 ***Quality control of purified protein libraries***

20 For mass spectrometry, the purified protein library sample was resuspended in water. The
21 spectrum was collected after addition of 2,5-dihydroxybenzoic acid matrix substance (Merck) using
22 an UltrafleXtreme™ MALDI-TOF/TOF mass spectrometer (Bruker Daltonics, Germany) in linear
23 mode.

24

25 **Acknowledgements**

26 We are grateful to Prof. Hideki Taguchi and Prof. Tatsuya Niwa for kindly providing us with the
27 expression plasmid of the Lon protease used in this study. We would also like to acknowledge
28 Dan S. Tawfik[†] and Valerio Guido Giacobelli for helpful discussions regarding this manuscript. In
29 addition, we would like to thank Kateřina Nováková for her technical help with collecting MALDI
30 spectra. This work was supported by the Czech Science Foundation (GAČR) grant number 17-
31 10438Y and the Human Frontier Science Program grant HFSP-RGY0074/2019. K.F. is supported
32 by ELSI - First Logic Astrobiology Donation Program.

33 **Competing Interest Statement:** The authors declare no competing interests.
34

35 **Author Contributions:** VT and KH designed research; VT, JVy, TN and KF performed
36 research; VT, JVo, KF and KH analyzed data; VT and KH wrote the paper.

37

38

39

1 References

- 2 1. Cleaves, H. J. The origin of the biologically coded amino acids. *J. Theor. Biol.* **263**, 490–
3 498 (2010).
- 4 2. Zaia, D. A. M., Zaia, C. T. B. V. & De Santana, H. Which amino acids should be used in
5 prebiotic chemistry studies? *Orig. Life Evol. Biosph.* **38**, 469–488 (2008).
- 6 3. Higgs, P. G. & Pudritz, R. E. A thermodynamic basis for prebiotic amino acid synthesis
7 and the nature of the first genetic code. *Astrobiology* **9**, 483–490 (2009).
- 8 4. Trifonov, E. N. Consensus temporal order of amino acids and evolution of the triplet code.
9 *Gene* **261**, 139–151 (2000).
- 10 5. Weber, A. L. & Miller, S. L. Reasons for the occurrence of the twenty coded protein amino
11 acids. *J. Mol. Evol.* **17**, 273–284 (1981).
- 12 6. Freeland, S. 'Terrestrial' Amino Acids and their Evolution. *Amin. Acids, Pept. Proteins*
13 *Org. Chem.* **1**, 43–75 (2010).
- 14 7. Philip, G. K. & Freeland, S. J. Did evolution select a nonrandom 'alphabet' of amino
15 acids? *Astrobiology* **11**, 235–240 (2011).
- 16 8. Ilardo, M., Bose, R., Meringer, M., Rasulev, B., Grefenstette, N., Stephenson, J.,
17 Freeland, S., Gillams, R. J., Butch, C. J. & Cleaves, H. J. Adaptive Properties of the
18 Genetically Encoded Amino Acid Alphabet Are Inherited from Its Subsets. *Sci. Rep.* **9**,
19 (2019).
- 20 9. Pace, N. R. The universal nature of biochemistry. *Proc. Natl. Acad. Sci. U. S. A.* **98**, 805–
21 808 (2001).
- 22 10. Holliday, G. L., Fischer, J. D., Mitchell, J. B. O. & Thornton, J. M. Characterizing the
23 complexity of enzymes on the basis of their mechanisms and structures with a bio-
24 computational analysis. *FEBS J.* **278**, 3835–3845 (2011).
- 25 11. Di Mauro, E., Dunker, A. K. & Trifonov, E. N. Disorder to Order, Nonlife to Life: In the
26 Beginning There Was a Mistake. 415–435 (2012).
- 27 12. Newton, M. S., Morrone, D. J., Lee, K. H. & Seelig, B. Genetic Code Evolution
28 Investigated through the Synthesis and Characterisation of Proteins from Reduced-
29 Alphabet Libraries. *ChemBioChem* **20**, 846–856 (2019).
- 30 13. Tanaka, J., Doi, N., Takashima, H. & Yanagawa, H. Comparative characterization of
31 random-sequence proteins consisting of 5, 12, and 20 kinds of amino acids. *Protein Sci.*
32 **19**, 786–795 (2010).
- 33 14. Riddle, D. S., Santiago, J. V., Bray-Hall, S. T., Doshi, N., Grantcharova, V. P., Yi, Q. &
34 Baker, D. Functional rapidly folding proteins from simplified amino acid sequences. *Nat.*
35 *Struct. Biol.* **4**, 805–809 (1997).
- 36 15. Longo, L. M., Lee, J. & Blaber, M. Simplified protein design biased for prebiotic amino

- 1 acids yields a foldable, halophilic protein. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 2135–2139
2 (2013).
- 3 16. Shibue, R., Sasamoto, T., Shimada, M., Zhang, B., Yamagishi, A. & Akanuma, S.
4 Comprehensive reduction of amino acid set in a protein suggests the importance of
5 prebiotic amino acids for stable proteins. *Sci. Rep.* **8**, (2018).
- 6 17. Solis, A. D. Reduced alphabet of prebiotic amino acids optimally encodes the
7 conformational space of diverse extant protein folds. *BMC Evol. Biol.* **19**, 1–19 (2019).
- 8 18. Giacobelli, V., Fujishima, K., Lepšík, M., Tretyachenko, V., Kadavá, T., Bednárová, L.,
9 Novák, P. & Hlouchová, K. In vitro evolution reveals primordial RNA-protein interaction
10 mediated by metal cations. *bioRxiv* (2021).
- 11 19. Longo, L. M., Despotović, D., Weil-Ktorza, O., Walker, M. J., Jabłońska, J., Fridmann-
12 Sirkis, Y., Varani, G., Metanis, N. & Tawfik, D. S. Primordial emergence of a nucleic acid-
13 binding protein via phase separation and statistical ornithine-to-arginine conversion. *Proc.*
14 *Natl. Acad. Sci. U. S. A.* **117**, 15731–15739 (2020).
- 15 20. Houben, B., Michiels, E., Ramakers, M., Konstantoulea, K., Louros, N., Verniers, J., der
16 Kant, R., De Vleeschouwer, M., Chicória, N., Vanpoucke, T., Gallardo, R., Schymkowitz,
17 J. & Rousseau, F. Autonomous aggregation suppression by acidic residues explains why
18 chaperones favour basic residues. *EMBO J.* **39**, 1–22 (2020).
- 19 21. Chiarabelli, C., Vrijbloed, J. W., Thomas, R. M. & Luisi, P. L. Investigation of de novo
20 Totally Random Biosequences. *Chem. Biodivers.* **3**, 827–839 (2006).
- 21 22. Keefe, A. D. & Szostak, J. W. Functional proteins from a random-sequence library.
22 *Nature* **410**, 715–718 (2001).
- 23 23. Labean, T. H., Butt, T. R., Kauffman, S. A. & Schultes, E. A. Protein folding absent
24 selection. *Genes (Basel)*. **2**, 608–26 (2011).
- 25 24. Yu, J. F., Cao, Z., Yang, Y., Wang, C. L., Su, Z. D., Zhao, Y. W., Wang, J. H. & Zhou, Y.
26 Natural protein sequences are more intrinsically disordered than random sequences.
27 *Cell. Mol. Life Sci.* **73**, 2949–2957 (2016).
- 28 25. Tretyachenko, V., Vymětal, J., Bednárová, L., Kopecký, V., Hofbauerová, K., Jindrová,
29 H., Hubálek, M., Souček, R., Konvalinka, J., Vondrášek, J. & Hlouchová, K. Random
30 protein sequences can form defined secondary structures and are well-tolerated in vivo.
31 *Sci. Rep.* **7**, (2017).
- 32 26. Tretyachenko, V., Voracek, V., Soucek, R., Fujishima, K. & Hlouchova, K. CoLide:
33 Combinatorial library design tool for probing protein sequence space. *Bioinformatics* **37**,
34 482–489 (2021).
- 35 27. Bateman, A. *et al.* UniProt: The universal protein knowledgebase in 2021. *Nucleic Acids*
36 *Res.* **49**, D480–D489 (2021).
- 37 28. Vymětal, J., Vondrášek, J. & Hlouchová, K. Sequence versus composition: What
38 prescribes IDP biophysical properties? *Entropy* **21**, 1–8 (2019).

- 1 29. Johannes, S., Ignacio, A.-C., Erwin, F., Verena, K., Mark, L., Tobias, P., Stephan, P.,
2 Curtis, R., Stephan, S., Benjamin, S., Jean-Yves, T., Daniel James, W., Volker, H., Kevin,
3 E., Pavel, T. & Albert, C. Fiji: an open-source platform for biological-image analysis. *Nat.*
4 *Methods* **9**, 676–682 (2012).
- 5 30. Melderer, L. Van & Aertsen, A. Regulation and quality control by Lon-dependent
6 proteolysis. *Res. Microbiol.* **160**, 645–651 (2009).
- 7 31. Niwa, T., Uemura, E., Matsuno, Y. & Taguchi, H. Translation-coupled protein folding
8 assay using a protease to monitor the folding status. *Protein Sci.* **28**, 1252–1261 (2019).
- 9 32. Bornberg-Bauer, E., Hlouchova, K. & Lange, A. Structure and function of naturally
10 evolved de novo proteins. *Curr. Opin. Struct. Biol.* **68**, 175–183 (2021).
- 11 33. White, S. H. The evolution of proteins from random amino acid sequences: II. Evidence
12 from the statistical distributions of the lengths of modern protein sequences. *J. Mol. Evol.*
13 **38**, 383–394 (1994).
- 14 34. Neme, R., Amador, C., Yildirim, B., McConnell, E. & Tautz, D. Random sequences are an
15 abundant source of bioactive RNAs or peptides. *Nat. Ecol. Evol.* **1**, 1–7 (2017).
- 16 35. Bornberg-Bauer, E., Hlouchova, K. & Lange, A. Structure and function of naturally
17 evolved de novo proteins. *Curr. Opin. Struct. Biol.* **68**, 175–183 (2021).
- 18 36. Davidson, A. R. & Sauer, R. T. Folded proteins occur frequently in libraries of random
19 amino acid sequences. *Proc. Natl. Acad. Sci. U. S. A.* **91**, 2146–2150 (1994).
- 20 37. Benner, S. A. Enzyme Kinetics and Molecular Evolution. *Chem. Rev.* **89**, 789–806
21 (1989).
- 22 38. Raggi, L., Bada, J. L. & Lazcano, A. On the lack of evolutionary continuity between
23 prebiotic peptides and extant enzymes. *Phys. Chem. Chem. Phys.* **18**, 20028–20032
24 (2016).
- 25 39. Brack, a & Orgel, L. E. Beta structures of alternating polypeptides and their possible
26 prebiotic significance. *Nature* **256**, 383–387 (1975).
- 27 40. Lupas, A. N. & Alva, V. Ribosomal proteins as documents of the transition from
28 unstructured (poly)peptides to folded proteins. *J. Struct. Biol.* **198**, 74–81 (2017).
- 29 41. Kovacs, N. A., Petrov, A. S., Lanier, K. A. & Williams, L. D. Frozen in Time: The History of
30 Proteins. *Mol. Biol. Evol.* **34**, 1252–1260 (2017).
- 31 42. Makarov, M., Meng, J., Tretyachenko, V., Srb, P., Březinová, A., Giacobelli, V. G.,
32 Bednárová, L., Vondrášek, J., Dunker, A. K. & Hlouchová, K. Enzyme catalysis prior to
33 aromatic residues: Reverse engineering of a dephosphoCoA kinase. *bioRxiv* (2020).
- 34 43. Kimura, M. & Akanuma, S. Reconstruction and Characterization of Thermally Stable and
35 Catalytically Active Proteins Comprising an Alphabet of ~ 13 Amino Acids. *J. Mol. Evol.*
36 **88**, 372–381 (2020).

- 1 44. Longo, L. M., Tenorio, C. A., Kumru, O. S., Middaugh, C. R. & Blaber, M. A single
2 aromatic core mutation converts a designed 'primitive' protein from halophile to
3 mesophile folding. *Protein Sci.* **24**, 27–37 (2015).
- 4 45. Despotović, D., Longo, L. M., Aharon, E., Kahana, A., Scherf, T., Gruic-Sovulj, I. &
5 Tawfik, D. S. Polyamines mediate folding of primordial hyperacidic helical proteins.
6 *Biochemistry* **59**, 4456–4462 (2020).
- 7 46. Yadid, I., Kirshenbaum, N., Sharon, M., Dym, O. & Tawfik, D. S. Metamorphic proteins
8 mediate evolutionary transitions of structure. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 7287–
9 7292 (2010).
- 10 47. Islas, S., Velasco, A. M., Becerra, A., Delaye, L. & Lazcano, A. Hyperthermophily and the
11 origin and earliest evolution of life. *Int. Microbiol.* **6**, 87–94 (2003).
- 12 48. Heffernan, R., Paliwal, K., Lyons, J., Singh, J., Yang, Y. & Zhou, Y. Single-sequence-
13 based prediction of protein secondary structures and solvent accessibility by deep whole-
14 sequence learning. *J. Comput. Chem.* **39**, 2210–2216 (2018).
- 15 49. Jones, T. Protein secondary structure prediction based on position-specific scoring
16 matrices. *J Mol Biol* **292**, (1999).
- 17 50. Frishman, D. & Argos, P. Seventy-five percent accuracy in protein secondary structure
18 prediction. *Proteins Struct. Funct. Genet.* **27**, 329–335 (1997).
- 19 51. Cuff, J. a & Barton, G. J. Application of multiple sequence alignment profiles to improve
20 protein secondary structure prediction. *Proteins* **40**, 502–511 (2000).
- 21 52. Levin, J. M. Exploring the limits of nearest neighbour secondary structure prediction.
22 *Protein Eng.* **10**, 771–776 (1997).
- 23 53. Garnier, J., Gibrat, J. F. & Robson, B. GOR method for predicting protein secondary
24 structure from amino acid sequence. *Methods Enzymol.* **266**, 540–553 (1996).
- 25 54. Fang, Y., Gao, S., Tai, D., Middaugh, C. R. & Fang, J. Identification of properties
26 important to protein aggregation using feature selection. *BMC Bioinformatics* **14**, 314
27 (2013).
- 28 55. Afgan, E. *et al.* The Galaxy platform for accessible, reproducible and collaborative
29 biomedical analyses: 2018 update. *Nucleic Acids Res.* **46**, W537–W544 (2018).
- 30 56. Rebollo, I. R., Sabisz, M., Baeriswyl, V. & Heinis, C. Identification of target-binding
31 peptide motifs by high-throughput sequencing of phage-selected peptides. *Nucleic Acids*
32 *Res.* **42**, e169–e169 (2014).

33