

Evolution of Networks of Protein Domain Organization

M. Fayez Aziz

University of Illinois at Urbana Champaign

Gustavo Caetano-Anollés (✉ gca@illinois.edu)

University of Illinois at Urbana Champaign

Research Article

Keywords: Domains, multidomains, evolution, time events, age, network, connectivity, modularity, randomness, scale-free, scale-rich

Posted Date: December 8th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-119891/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 Evolution of Networks of Protein Domain Organization

2 M. Fayez Aziz and Gustavo Caetano-Anollés*

3 Evolutionary Bioinformatics Laboratory, Department of Crop Sciences, University of Illinois,
4 Urbana, IL 61801, United States

5 *Correspondence: gca@illinois.edu (GCA)

6

7 Abstract (196 words)

8 Domains are the structural, functional and evolutionary units of proteins. They combine
9 to form multidomain proteins. The evolutionary history of this molecular combinatorics
10 has been studied with phylogenomic methods. Here, we construct networks of domain
11 organization and explore their evolution. These networks revealed two ancient waves of
12 structural novelty arising from ancient 'p-loop' and 'winged helix' domains and a massive
13 'big bang' of domain organization. The evolutionary recruitment of domains was highly
14 modular, hierarchical and ongoing. Domain rearrangements elicited non-random and
15 scale-free network structure. Comparative analyses of preferential attachment,
16 randomness and modularity of networks showed yin-and-yang complementary transition
17 patterns along the evolutionary timeline. Remarkably, evolving networks highlighted a
18 central evolutionary role of cofactor-supporting structures of non-ribosomal peptide
19 synthesis (NRPS) pathways, likely crucial to the early development of the genetic code.
20 Some highly modular domains featured dual response regulation in two-component
21 signal transduction systems with DNA-binding activity linked to transcriptional
22 regulation of responses to environmental change. Interestingly, hub domains across the
23 evolving networks shared the historical role of DNA binding and editing, an ancient
24 protein function in molecular evolution. Our investigation unfolds historical source-sink
25 patterns of evolutionary recruitment that further our understanding of protein
26 architectures and functions.

27

28 **Index terms:** Domains, multidomains, evolution, time events, age, network, connectivity,
29 modularity, randomness, scale-free, scale-rich.

30 Introduction

31 The biological functions of genes manifest through the proteins or functional RNA
32 molecules they encode. In evolution, novel functions appear when genes produce new
33 genes by duplication, mutation, recombination, fusion and fission, or when genes are
34 generated *de novo*. Research has attempted to quantitatively describe the origins of these
35 processes of molecular diversification and how they increase molecular complexity over
36 the course of evolution, for instance through pathways of protein domain organization^{1,2}.
37 Other examples include uncovering the natural history of biocatalysis by tracing chemical
38 mechanisms in enzymatic reactions³, evolutionary analysis of optimization and increase
39 of protein folding speed derived from a flexibility-correlated factor known as contact

40 order (the average relative distance of amino acid contacts in the tertiary structure of
41 proteins)⁴, and the study of biphasic-rewiring and modularity of metabolomic networks
42 of *Escherichia coli* minutes after subjection to stress⁵. In particular, the history of an
43 ‘elementary functionome’ was traced with a bipartite network of elementary functional
44 loop sequences and structural domains of proteins⁶. The study revealed two initial waves
45 of functional innovation involving founder ‘p-loop’ and ‘winged helix’ domain structures
46 and the emergence of hierarchical modularity and power law behavior in network
47 evolution.

48 Protein domains are structural and functional units of evolution that make up proteins⁷,
49 sometimes in unusually complex arrangements⁸. They fold into compact 3-dimensional
50 (3D) atomic structures that arrange alpha-helical and beta-sheet structure elements into
51 tightly packed conformations of the polypeptide chain. The Structural Classification of
52 Proteins (SCOP)⁹ and its extended version SCOPe¹⁰ are popular taxonomy gold standards
53 of domain structure. SCOP definitions can be used to scan genome sequences for motifs
54 of domains and study how they combine in proteins⁷. In SCOP, the structure of domains
55 exhibiting similar 3D arrangements of secondary structures and thus identical topologies
56 have been classified as folds (F)⁹. Within folds, protein domains whose structure and
57 functional features indicate a common evolutionary origin are further grouped into fold
58 superfamilies (FSF). These FSFs sometimes hold multiple evolutionarily related families
59 (Supplementary Fig. S1A). As of June 3, 2020, 276,231 annotated SCOPe domains
60 populate the 164,840 protein structures of the Research Collaboratory for Structural
61 Bioinformatics Protein Data Bank (RCSB-PDB).

62 Domain structures appear repeatedly in the protein molecules, singly or in combination
63 with other domains⁸. More than two-thirds of protein sequences are longer than an
64 average domain length, a vast majority of which are multidomain proteins¹¹. A study of
65 protein structures in 745 genomes showed that the lengths of orthologous protein
66 families in Eukarya were almost double the lengths found in Bacteria and Archaea¹². This
67 variance among lengths results from shorter prokaryotic nondomain sequences that link
68 domains to each other in proteins and have evolved reductively in prokaryotes but not in
69 eukaryotes. The arrangement of domains along the sequence of multimeric proteins is
70 referred to as ‘*domain organization*.’ Both the structure and organization of domains,
71 which have been collectively termed protein domain ‘*architecture*’, are considered far
72 more evolutionarily conserved than protein sequence^{8,13,14}. In addition, some domain
73 combinations make up functional units that recur in different protein contexts¹⁵. They
74 have been termed supradomains (Supplementary Fig. S1B). Thus, domains and domain
75 combinations behave as modules, parts that interact with each other more than with
76 other parts or modules of the system.

77 The evolution of protein domain architecture can be studied with phylogenomic
78 methods^{16,17}. One approach takes advantage of the reconstruction of phylogenomic trees
79 from the occurrence and abundance of architectures in proteomes at F and FSF levels
80 based on the sequence and structure of millions of protein sequences encoded in
81 hundreds of genomes^{18,19}. The trees have leaves representing single domain and

82 multidomain proteins. They allow to build evolutionary timelines of molecular accretion
83 using a phylogenomic framework that describes the evolutionary history of a growing
84 molecular interactome⁸. Remarkably, studies showed that architectural diversification
85 evolved through gradual accumulation of domains (singly occurring domains), domain
86 pairs (two different domains), multidomains (numerous domains, with occasional
87 repetition) and domain repeats (domains of one type that are repeated)⁸. The
88 diversification began with a few single-domain architectures earlier in the timeline,
89 followed by an increasing rate of accretion that culminated in a massive “big bang” of
90 domain combinations. The accumulation of architectures continued to date but with a
91 decreasing rate^{6,8}.

92 Here, we continue to explore the evolving interactome of protein domain organization.
93 We use the phylogenomic tree of architectures⁸ to generate a timeline that captures the
94 historical development of domain and multidomain interactions with a graph theoretical
95 approach⁷ of evolving network structure. The timeline was calibrated with a molecular
96 clock of protein structures, which assigns relative ages of domains to billions of years (Gy)
97 of geological time²⁰. Five distinct composition- and topology-based ‘operative’ criteria of
98 connectivity defined nodes and links of the evolving networks. This strategy identified
99 connectivity distributions in a series of 169 growing networks, hubs of evolutionary
100 recruitment acting as *donors* and *acceptors*, and structural adaptations of evolving
101 networks to modular, random and scale-free properties. In particular, we discover a
102 pattern of connectivity driven by fusions and fissions, respectively, with densely linked
103 older and younger architectures from the evolutionary timeline sandwiching a period of
104 sparse connectivity. This supports a biphasic or hourglass pattern previously observed in
105 protein evolution²¹ and follows a model of module emergence²². We thus reveal
106 remarkable patterns of emergence of hierarchy, modularity and structural cooption in
107 evolving networks.

108 **Results and discussion**

109 *Construction of evolving networks*

110 We build evolving networks of domain organization to explore how single-domain and
111 multidomain proteins share domain make-up and how recruitment processes shape
112 protein evolution. An ‘entity set’ of domains, supradomains, and multidomains were first
113 extracted from the genomic census of fold structure and domain organization. This set of
114 component parts of proteins, mostly recurrent, defined the nodes of the networks, which
115 were labeled with *concise classification strings (ccs)* describing SCOP domain constituents
116 (Fig. 1A). We define supradomains as sub-combinations of domains that appear in the
117 census and are often used as evolutionary building blocks of multidomains. The
118 definition is more inclusive than that of ref¹⁵.

119 The growing interactions among contemporary architectures were captured with five
120 different operative criteria for network generation defined by composition, pairwise
121 occurrence, adjacency, and splicing of domain parts in a protein molecule, where: (i)

122 *composition* describes makeup (component parts) of the molecular whole; (ii) *pairwise*
123 *occurrence* describes appearance of parts in sets of two; (iii) *adjacency* refers to their
124 geometrical or spatial arrangement (topology); and (iv) *splicing* refers to the
125 rearrangement of parts by operations of joining and excision that decompose structures
126 (Fig. 1B). The Composition Network (CX) linked domain and supradomain to
127 multidomain nodes (in a partially bimodal fashion) when proteins shared compositional
128 makeup. The Pairwise Network (PX) connected domain to supradomain nodes when
129 components occurred in pairs in a protein. The Pairwise Adjacency Network (PAX)
130 connected domain to supradomain nodes when components occurred in pairs that were
131 adjacent. The Spliced Pairwise Network (SPX) linked domain nodes to each other when
132 their pairs were present in domain-spliced proteins. Lastly, the Spliced Pairwise
133 Adjacency Network (SPAX) linked domain nodes to each other when their adjacent pairs
134 were present in the domain-spliced proteins (Fig. 2).

135 We then mapped the evolutionary ages of architectures onto the nodes of networks built
136 using these five operative criteria (Supplementary Fig. S2). We did so for each of the 169
137 time-events of the timeline. Network construction has been illustrated with connectivity
138 details of the most ancient domains (Supplementary Fig. S3) and further described in
139 Section 1 of Supplementary Text. Networks showcased time directionality, connectivity
140 distributions, and network layouts:

141 (i) *Time Directionality*: Mapping ages onto networks helped follow their evolutionary
142 growth, as nodes and links accumulated over time since the origin of proteins to the
143 present. The timeline of networks imposed a time directionality on network links, making
144 them *arcs* (directed edges with arrows pointing from older to younger nodes) of directed
145 graphs (Fig. 1C). The ages of arcs were borrowed from the youngest of the component
146 nodes involved in a link (Supplementary Fig. S3B).

147 (ii) *Degree Distributions*: The number of links connected to a node define that node's
148 'degree'. The degree distribution is a 'composability' attribute of a network and the entity
149 set represented by its nodes, a design principle describing the inter-relationship of
150 components of a system. In network evolution, the appearance of a new node may trigger
151 establishment of one or more arcs from existing (older) nodes. Furthermore, *outdegree*
152 describes the number of outward links and *indegree* the number of inward links from a
153 node. As the timeline progresses, older nodes gain higher outdegrees as compared to the
154 higher indegrees of recent nodes (Fig. 1C), polarizing the network with arcs depicting
155 'arrows of time' (Supplementary Figs. S2 and S3). The chronological appearance of
156 architectures (domains, supradomains and multidomains) as network connectivity
157 expands along the timeline causes degree to accumulate in the evolving networks (Fig. 2).
158 Multiple interactions of nodes along the timeline diversified connectivity, a feature
159 captured and quantified by weighted degree. Interestingly, box-and-whisker's plots of
160 weighted outdegree and indegree demonstrate bimodal degree distributions typical of
161 biological systems^{5,22} (Supplementary Fig. S4). The yin-yang patterns of contractions and
162 expansions of architectural innovation are evident from the distributions of modern
163 outdegrees and indegrees (Supplementary Fig. S5). In particular, the cumulative

164 outdegree and indegree scattergrams demonstrate an hourglass (or bimodal) pattern of
165 linkage development unfolding in evolution (Supplementary Fig. S6).

166 (iii) *Time Event-based ‘Radial’ and ‘Waterfall’ Layouts*: The growth of a network evolving
167 at discrete temporal intervals can be modeled with Discrete Event Simulation (DES)
168 tools^{23–25}. Borrowing the DES rationale, we modeled the evolution of directed networks of
169 domain organization with time flowing from one event to another as discrete
170 evolutionary ‘time steps’, typical of a *step function*. The progression of events was
171 visualized with two types of layouts, a vertical representation we coined ‘*waterfall*’ layout
172 that had nodes arranged top-down by age and a concentric ‘*radial*’ representation of
173 growing networks that unfolded time-events of protein evolution from center to
174 periphery (Fig. 1C). Network clusters comprising of hubs and their cohesive neighbors
175 were segregated to improve differentiation along the horizontal axis. The *waterfall* and
176 *radial* layouts made evolutionary recruitment evident as time events progressed
177 downward or outward, respectively (Figs. 2 and 3).

178 *Early history of modern domain organization*

179 The accumulation of links connecting domain, supradomain and multidomain proteins
180 in evolving CX, PX, PAX, SPX and SPAX networks played back the complicated history
181 of domain recruitments that drive the evolution of domain organization. Figure 2 shows
182 networks in radial layout at representative time-events defining boundaries of the three
183 epochs of the evolving protein world (‘architectural diversification’, ‘superkingdom
184 specification’ and ‘organismal diversification’, *sensu*^{8,19}). Networks grew in time and
185 became increasingly complicated tangles, massively expanding after a “big bang” of
186 domain combinations during the organismal diversification epoch. Movies described the
187 evolutionary growth of these networks (Supplementary Video 1).

188 To illustrate the versatility of the waterfall visualization strategy, we dissected the early
189 origin of proteins with the SPX network. Two major waves of structural innovation
190 arising from ancient ‘*p*-loop’ and ‘winged helix’ domains were observed in the waterfall
191 diagrams of a highly connected (reduced) subnetwork visualization of the SPX network
192 (Fig. 3), matching similar recruitment waves observed in the study of evolutionary
193 networks of elementary functionomes⁶ and metabolites²⁶. Waves originated in primordial
194 $\alpha/\beta/\alpha$ -layered sandwich, β -barrel and helical bundle structures identified in an earlier
195 structural phylogenomic study as part of the most ancient 54 protein domain families²⁷.
196 However, most of the connectivity of these major pathways was established during the
197 organismal diversification epoch less than 1.5 Gy ago ($nd \geq 0.6$) and hence was fully
198 developed relatively recently in evolution. The ‘*p*-loop’ and ‘winged helix’ waves
199 embedded the major gateways of enzymatic recruitment we previously reported for
200 metabolism²⁶. The first gateway was mediated by the c.37 P-loop hydrolase fold and
201 originated in the energy interconversion pathways of the purine metabolism subnetwork.
202 The second pathway was mediated by the a.4 winged helix fold and originated in the
203 biosynthesis of cofactors and the metabolic subnetwork of porphyrin and
204 chlorophyll^{16,26,28}. The congruent realization of these evolutionary patterns with data

205 sources of different types is remarkable (Supplementary Video 2). It strongly supports the
206 historical statements we propose. Further information can be found in Section 2 of
207 Supplementary Text.

208 *Network analysis of cooption mechanisms of recruitment*

209 The networks of domains (SPX and SPAX) elicited 161 unique time-events along the
210 evolutionary timeline, out of a total 169 events expected for networks of domains,
211 supradomains and domain combinations (CX, PX and PAX) (Supplementary Tables 1-5).
212 The node and connectivity distributions among the time-event bins of the evolving
213 networks highlight the widespread, growing and recurrent combinatorial recruitment
214 process that incorporates domains and their combinations into protein scaffolds and
215 drives structural evolution (Fig. 2). Indeed, the largest hubs representing the most
216 popular domains in the highly connected SPX subnetwork appeared not only early in
217 evolution but also in the modern protein world (Fig. 3). Similar to the evolution of
218 elementary functions⁶, domain innovation also developed early during the first ~1.8 Gy of
219 protein history (Fig. 3). The combinatorial recruitment process however spanned the
220 entire timeline (Supplementary Fig. S2). In terms of origins, the first donor and acceptor
221 composition event occurred in protein evolution with the appearance of a link in the CX
222 network connecting domain c.2.1 to domain combination c.2.1|a.100.1, ~3.54 Gya ($nd =$
223 0.069). The first donor and acceptor pair occurred in the pairwise PX and SPX networks
224 ~3.12 Gya ($nd = 0.179$), ~0.42 Gy later ($\Delta nd = 0.11$). The pairing event involved domains
225 c.37.1 and d.14.1. The first adjacent donor and acceptor pair of the adjacency-based PAX
226 and SPAX networks appeared ~2.90 Gya ($nd = 0.237$), ~0.22 Gy later ($\Delta nd = 0.06$). The
227 adjacently paired nodes were domains c.37.1 and c.23.16. These observations highlight a
228 remarkable tendency of domain organization to gradually but recurrently constrain
229 pairwise occurrences in multidomain proteins. The evolutionary history of donors and
230 acceptors of domain organization is hence associated with a highly optimized process of
231 cooption. To explore this combinatorics, first we dissected the network connectivity with
232 bar plots that describe the chronological accumulation of links along the evolutionary
233 timeline (Supplementary Fig. S7). This made general patterns quantitative and source-
234 sink relationships explicit. Second, we analyzed the per unit donor/acceptor ratio in the
235 evolving networks to highlight pairwise cooption and composability, respectively
236 (Supplementary Fig. S8). Specifically, domain acceptors (represented by network
237 indegree) of SPX increased in number to a global average of 8.63 (± 0.15) sinks per
238 domain in evolution. Domain donors (represented by network outdegree) of SPX reached
239 a higher global average of 9.7 (± 0.56) sources per domain, indicating significant
240 reutilization of relatively ancient domains. In contrast, the average number of donors and
241 acceptors in the evolving CX network plateaued at 3.41 ± 0.34 sources and 3.43 ± 0.05 sinks
242 per domain/multidomain, respectively. This showed uniform source/sink evolutionary
243 rates as proteins acquired higher composability with time. Third, an inferential analysis of
244 cooption-based source-sink relationships maturing at modern times revealed an
245 independence of patterns from the selected network generation criteria (Supplementary
246 Fig. S9). Primarily, the composition events yielding source domains and supradomains
247 were dominant, with the number of events almost doubling in the CX network from the

248 origin to the organismal diversification epoch ~ 1.5 Gya ($nd = 0.6$). However, the pairwise
249 cooption events of the SPX domain network, e.g., doubled in number and reached
250 relatively comparable levels in evolution only after delays of ~ 0.6 Gy ($\Delta nd = 0.15$) and
251 ~ 2.1 Gya ($nd = 0.75$), respectively. Moreover, the number of cooption events yielding sink
252 domains in SPX almost tripled by the beginning of the organismal diversification epoch.
253 In contrast, the number of CX sinks reached that level only halfway along that
254 evolutionary epoch. These divergent patterns indicate a frustrated dynamics of network
255 growth. The early adoption of composability of domains and supradomains in
256 multidomains seems to have preceded the pairwise cooption of domains in protein
257 history, leading to the numerous recruitment pathways of the modern protein world. A
258 discussion on the source-sink relationships impacted by domain fusion and fission
259 processes can be found in *Section 3* of [Supplementary Text](#).

260 *Hubs in network evolution*

261 Network hubs are at the heart of network connectivity and could chaperone network
262 evolution²⁹. We ranked modern domains and domain combinations of age $nd = 1$ as hubs
263 based on the 99.9th percentile of indegree and outdegree. Hubs were annotated with
264 domain organization attributes, including SCOP domain descriptions, age,
265 fusional/fissional information, and GO terms. We also associated hubs with age ranks
266 reflecting their order of evolutionary appearance in the timeline.

267 The most notable donor hubs for all networks types were the carrier protein domains
268 e.23.1, a.28.1 and c.69.1, which are involved in Non-Ribosomal Peptide Synthesis (NRPS),
269 whether directly or indirectly through other pathways ([Table 1](#)). These domains
270 diversified later in evolution yielding cofactor-binding molecular switches and barrel
271 structures²⁷. Ancient NRPS pathways of domain accretion have been associated with a
272 model that not only described stabilization and decoration of membranes by primordial
273 alpha-helical bundles and beta-sheets, but also explained primordial protein synthesis
274 and genetic code specificity chaperoned by ancient forms of aminoacyl-tRNA synthetase
275 (aaRS) catalytic domains and NRPS modules. NRPS even preceded the emergence of the
276 ribosome, acting as scaffold for nucleic acids and the modern translation function. In
277 particular, the PX and PAX networks highlight the central evolutionary role of these
278 novel emerging cofactor structures in the NRPS pathways. Thus, our findings made
279 explicit that our connectivity criteria of generating networks of domain organization were
280 at the cornerstone of the early development of genetic code and supported the
281 evolutionary model of early biochemistry based on phylogenomic information and
282 network structure.

283 Domains c.30.1, b.1.1, d.142.1 and g.3.11 ($0.723 < nd < 0.977$) were the most prominent
284 acceptor hubs ([Table 2](#)). These structures are integral parts of two-component signal
285 transduction systems that are common in microbes. The highly modular domains feature
286 dual response regulator proteins involved in the two-component signal transduction
287 system comprising of an N-terminal response regulator receiver domain and a variable C-
288 terminal effector domain with DNA-binding activity. These proteins are transcriptional

289 regulators in bacteria and some protozoa, detecting and responding to environmental
 290 changes, e.g. nitrogen fixation. These evolving interactions of microbes with the
 291 environment mediated by two-component systems have apparently influenced the
 292 evolutionary process of cooption. Three acceptor hubs that were significant in PX with
 293 indegree > 250 (following behind the 99.9th percentile in other networks) were Nucleotide
 294 cyclase (d.58.29), Spermadhesin, CUB domain (b.23.1), and Fibronectin type III (b.1.2)
 295 ($nd = 0.723-0.809$). See [Section 4 of Supplementary Text](#) for additional donor/acceptor
 296 hub information, and [Section 5](#) for cooption events occurring during the ‘big bang’ of
 297 domain organization.

298 *Emergence of preferential attachment in network evolution*

299 Genomic-centric processes such as duplication, recombination, fusion and fission shape
 300 patterns of molecular complexity². Many of these patterns can be explained with large
 301 ‘scale-free’ networks that grow by following the preferential attachment principle³⁰. These
 302 self-organizing and highly inhomogeneous networks attach links to highly connected
 303 hub-like nodes in a ‘rich-get-richer’ fashion, lacking a characteristic scale, irrespective of
 304 the properties of individual nodes or systems³¹. This pattern of network expansion, which
 305 is remarkably popular in biology³², is sharply distinct from that of the Erdős–Rényi
 306 random network model^{33,34}. In a scale-free network, the probability $P(k)$ of nodes
 307 connecting with neighboring k nodes (i.e. the ratio of nodes with k links) decays as a
 308 power law, $P(k) \sim k^{-\gamma}$, with γ defined as the exponent of power law decay. The frequency
 309 distributions of node-connectivity in biomolecular networks have γ typically ranging 2.1–
 310 2.4³⁵. Thus, scale-free properties drive degree distributions entailing heavy tails, where
 311 very few nodes have high degree values.

312 Our statistical analyses of the featured indegree distributions along the timeline of
 313 growing networks uncovered interesting patterns of power law dynamics ([Fig. 4](#)). The
 314 scale-free patterns were established early on in protein evolution, primarily evident in the
 315 CX composition network. These patterns were remarkably divergent from evolving
 316 networks connected at random (RVN p -value > 0.05). While power law behavior
 317 generally declined as the networks evolved (KS p -value < 0.05, α < 2.5), it somewhat
 318 sustained after the ‘big bang’ but only in CX and not in the pairwise networks (KS fit and
 319 γ closer to 0 and -2 in CX, respectively). A log linear regression model of CX produced
 320 the highest absolute value for γ of 3.81 among the five networks, which was achieved early
 321 along the evolutionary timeline ($nd \sim 0.25$). This value of γ was much higher than values
 322 reported for metabolic networks ($\gamma \sim 2.2$)³². Remarkably, the γ was maintained at ~ 3 before
 323 and after the ‘big bang’, while remaining at ~ 2 until modern times with a minimum value
 324 of 1.7. The other four networks generated primarily with pairwise criterion apparently
 325 deviated from the power-law behavior, especially after the ‘big bang’. For instance, the γ
 326 of PX and PAX peaked at 2.4 ($nd \sim 0.35$) and 3.2 ($nd \sim 0.38$), respectively, slightly later
 327 than CX. We also noted a transition in γ from 2.1 in PX and 2.7 in PAX prior to the ‘big
 328 bang’ to 1.6 in both after the big bang, plateauing at ~ 1 until the present. In the SPX and
 329 SPAX networks, γ reached a peak even later in time than PX and PAX with values of 2.8
 330 ($nd \sim 0.54$) and 3.4 ($nd \sim 0.66$), respectively. These values transitioned from 2.4 in SPX

331 and 2.8 in SPAX from before the big bang to 1.6 and 1.7 after the big bang, respectively,
 332 plateauing at ~ 1 in both the networks. As expected, the average γ based on less
 333 representative outdegree of each of the five networks remained low (1 ± 0.05).

334 We noticed that the patterns of γ curves over the connectivity of the networks were
 335 biphasic, with two minima at $nd \sim 0.37$ and ~ 0.67 . Moreover, the scale-free tendency of
 336 adjacency networks seemed comparatively higher than that of networks lacking the
 337 adjacency restriction. For instance, the average values of γ for the PAX and SPAX
 338 networks (1.87 ± 0.06 and 2.13 ± 0.07 , respectively) were relatively higher than those for
 339 the corresponding parent PX and SPX networks (1.61 ± 0.05 and 1.89 ± 0.06 ,
 340 respectively). This suggests that the proximity of residuals in the amino acid sequence
 341 plays a major role in rendering the power-law behavior of evolving networks of domain
 342 organization. Overall, the average γ of CX (2.56 ± 0.06) remained the highest along the
 343 evolutionary timeline, indicating that composition strongly elicits the preferential
 344 attachment property. A complementary transition from random to non-random behavior
 345 (RVN p-value: $1 \rightarrow 0$) in ancient networks ($nd \sim 0.3$) implies deviation from randomness
 346 as biological networks evolve. Remarkably, this transition event coincides with the origin
 347 of a processive ribosome. Such biphasic patterns are common in biology and have
 348 explained the emergence of biological modules²² in metabolic networks of *Escherichia*
 349 *coli*⁵, networks of elementary functionomes⁶, and molecular ancestry networks of
 350 enzymes³⁶. Section 6 of [Supplementary Text](#) further discusses *scale-freeness* and
 351 *randomness* of networks.

352

353 *Emergence of hierarchical modularity*

354

355 Modular networks embed sets of communities (closely-knit modules) that establish links
 356 preferentially within themselves and do so sparsely with the rest³⁷. Network modularity
 357 usually offsets the power-law behavior of biological networks by distributing node
 358 degrees within communities³⁸⁻⁴⁰. However, both scale-free properties and modular
 359 structure may co-exist in a network when modules coalesce hierarchically³². A primary
 360 index of modularity is the *average clustering coefficient* (C), defined as a node-averaged
 361 ratio of triangles (graph cycles of length 3) to triads (the connected graph triples) of the
 362 network, not taking into account the weights or direction of the node-links^{32,41,42} ([Fig. 5](#)).
 363 The adjacency PAX and SPAX networks both showed the lowest C (averaged over nd)
 364 with a value of 0.09 ± 0.009 . The composition CX network had a relatively higher C of 0.2
 365 ± 0.009 . However, the non-adjacency pairwise PX and SPX networks had the highest C
 366 values of 0.5 ± 0.02 and 0.32 ± 0.014 , respectively. These values were still lower than those
 367 reported for metabolic networks ($C = \sim 0.6$)^{32,40,43}. Hence, the networks supposedly
 368 evolved more random smaller modules connected by various inter-modular links, rather
 369 than stronger larger modules with few interconnections. Also, the evolution of modular
 370 structure appeared better consolidated by pairwise (PX and SPX) and to a lesser degree
 371 composability (CX) constraint rather than adjacency (PAX and SPAX) restriction.
 372 Comparing patterns of modularity of evolving networks to those of randomness (given by
 373 $RVN_{p\text{-value}}$) indicated complementary transitions between the two behaviors over the
 374 evolutionary timeline ([Figs. 4 and 5](#)).

375

376 In order to dissect the modular behavior of evolving networks, we studied the regression
377 patterns of C against network size N and evolutionary age nd . For typical scale-free
378 models, C declines sharply with increasing N ($C \sim N^{-\text{coefficient}}$), while the coefficients are as
379 high as 0.75⁴⁴. Instead, highly modular networks are typically independent of N ³². In our
380 networks, C regressed by N with very low coefficients (CX, 0.000036; PX, 0.00007; PAX,
381 0.000035; SPX, 0.00016; SPAX, 0.00016). In contrast, the regression of C with age ($C \sim nd^{\text{coefficient}}$)
382 produced significantly higher coefficients (CX, 0.39; PX, 0.85; PAX, 0.39; SPX,
383 0.35; SPAX, 0.41) (Fig. 5). The reference power-law (Barabási) networks that were used as
384 control showed a C of zero, as expected⁴⁵. Our data strongly suggests the existence of a
385 highly modular structure that is independent of network growth but is strongly
386 constrained by history, especially when considering the pairwise interactions of the PX
387 network. The rise of the modularity index with emerging power-law degree distribution
388 during certain periods of network evolution indicated a parallel formation of complex
389 hierarchical module clusters with scale-free properties, not distinct from those present in
390 metabolic networks³². Our networks of domain organization show a slight lag between an
391 onset of scale-free organization (measured with KS fit and γ indegree statistics) and a
392 delayed emergence of modular behavior (measured with C), occurring during early
393 protein evolution. This was followed by intermittent periods of hierarchical modularity
394 spanning across the middle of the evolutionary timeline. Remarkably, the evolving
395 networks showed a prominent biphasic pattern of hierarchical modularity involving two
396 peaks of modularity (higher statistic C) coinciding with increased power-law behavior
397 (valleys of KS fit and $-\gamma$ curves), at $nd \sim 0.37$ and $nd \sim 0.67$, respectively (Figs. 4 and 5).
398 The modularity heatmaps and dendrograms of select phases of network evolution
399 confirm these biphasic patterns (Fig. 6), which were markedly distinct from the long-
400 tailed clustering patterns of preferential attachment (Supplementary Fig. S10). As
401 identified earlier⁶, the timing of this switch coincides with the early development of
402 genetic code specificity in the emerging ribosomal aaRS catalytic domains, which was
403 facilitated by the OB-fold structure⁴⁶. These counteracting and delicately balanced trends
404 of modularity and preferential attachment suggest that the emergence of scale-free
405 behavior of the partial bipartite CX network must have impacted the hierarchical
406 modular structure of the modern pairwise networks of domain organization (PX, PAX,
407 SPX, SPAX) (Supplementary Video 3). A detailed account of our testing and verification
408 of this conjecture is explained in Section 7 of Supplementary Text.

409

410

Conclusions

411 We traced evolutionary ages inferred from a phylogenomic analysis of protein
412 architectures onto networks of domain organization. Evolving networks revealed two
413 prominent waves of structural novelty involving ancient domain innovations and founder
414 ‘ p -loop’ and ‘winged helix’ domain structures. We found that the evolutionary
415 recruitment of domains and multidomains in proteins was ongoing and highly modular.
416 Remarkably, the networks highlighted the role of cofactor-supporting structures of NRPS
417 pathways, which were backbone to the early evolution of the genetic code. The evolving
418 domain rearrangements featured multitier evolutionary episodes of scale-free network

419 structure, hierarchy and modular behavior. Remarkably, our analyses support biphasic
420 patterns of diversification and module emergence that we have observed earlier^{6,22}. In an
421 initial phase, at the cusp of architectural diversification, the modular components of
422 emerging domain organization associated through weak linkages of recruitment. The
423 second phase was massive and prolonged, with a multitude of modules appearing after
424 the ‘big bang’ of the protein world, supporting the onset of organismal diversification.
425 Such biphasic patterns are prevalent in biology and impact size, dipeptide makeup, and
426 loop-mediated flexibility of proteins, possibly due to their intrinsic disorder^{4,46}. Hence, the
427 existence of biphasic patterns in evolving networks might be integral to biological history.

428 **Methods**

429 *Experimental design*

430 *Phylogenomic analysis of the entity set of protein domain architectures*

431 We explore the evolution of networks describing how structural domains combine and
432 split to form single domain and multidomain proteins, i.e. the domain organization of
433 proteins. The definition of protein domain structures followed the FSF level of SCOP
434 version 1.75⁹ (Fig. 1). Domain interactions were studied along an evolutionary timeline of
435 structural and architectural innovation directly derived from a phylogenomic tree of
436 architectures reconstructed from a Hidden Markov-Model (HMM)-based census of
437 structural domain organization encoded in 1,730 FSF structures present in 749 genomes
438 of 52 archaeal, 478 bacterial and 219 eukaryal organisms (dataset A749)⁸ (Supplementary
439 Fig. S1). The phylogeny represents a reconstruction of the “natural history” of proteins
440 that is supported by a model of protein structural growth⁴⁷ and is carefully indexed with
441 various evolutionary epochs of the protein world⁸.

442 *Calculation of the ages of domain organization*

443 The ages of domains and domain combinations were calculated as node distance (*nd*)
444 values, which were derived directly from the rooted phylogenomic tree of protein domain
445 organization⁸. *nd* values describe relative ages (in a relative 0-1 scale) of first appearances
446 of 6,162 domains and domain combinations (multidomains) defined at SCOP FSF level
447 (the extant ‘entity set’ sampled by our study; Fig. 2) Collectively, ages defined an
448 evolutionary timeline embodying architectural transformations and molecular transitions
449 mediated by fusion and fission processes in the form of 169 unique ‘time events’ (age
450 groups or time slivers) (Supplementary Fig. S2). A Python script was used to count the
451 number of nodes from the root (base) of the tree to each leaf node and present the
452 distance matrix of nodes in a relative zero-to-one scale⁶. The script utilized the high
453 imbalance of phylogenomic trees as a fundamental feature to derive the relative ages of
454 domain organization⁸. The tree imbalance resulted from the accumulation of structures
455 and their combinations in proteins and proteomes and not from node density, thus
456 representing a true evolutionary process²⁰.

457 The timeline was calibrated with a molecular clock of FSF structures ($t = -3.831nd +$
458 3.628) used to calculate geological age in Gy through calibration points of FSF domains
459 associated with microfossil, fossil and biogeochemical evidence, biomarkers, and first-
460 appearance of clade-specific domains²⁰. The RSCB – PDB count was determined by
461 following the hyperlink associated to the number of entries or structures (which is
462 updated weekly) and selecting “Customizable Table” from the ‘Reports’ menu above the
463 results section. Subsequently, SCOP, CATH, and PFAM ID options were selected as
464 domain information under the ‘Domain Details’ section and domain counts data were
465 exported as a comma separated value (.csv) file report. [Supplementary Tables 1-5](#) provide
466 an exhaustive summary of various connectivity categories of evolving networks based on
467 this ‘entity set’ of domain organization. The extraction pipeline of SPX/SPAX domain
468 units from the original data set can be found in [Supplementary Table 6](#).

469 *Indexing domain attributes*

470 Domain ages and assignment of fusional/fissional properties followed ref. ⁸. SCOP *concise*
471 *classification strings (ccs)* of domain descriptions⁹ were downloaded from
472 <http://scop.mrc-lmb.cam.ac.uk/scop/parse/index.html> for SCOP version 1.75 as the file
473 `dir_des_scop_txt_1_75.txt`. Available descriptions for 2,223 single domains were obtained
474 from SCOP unique identifiers (sunID). The Gene Ontology (GO) specifications were
475 recorded from the Superfamily Database (SUPFAM) available at
476 <http://supfam.cs.bris.ac.uk/SUPERFAMILY/GO.html>. High-coverage domain-centric
477 GO annotations that were supported only by all UniProts (including multidomain
478 UniProts) were downloaded as the file `Domain2GO_supported_only_by_all.txt`. High-
479 quality truly domain-centric GO annotations that were supported by both single domain
480 UniProts and all UniProts (including multidomain UniProts) were downloaded as the file
481 `Domain2GO_supported_by_both.txt`. We reported only the GO annotations ‘by all’ to
482 capture higher coverage. Also, the GO terms were reported only for the 2,223 single
483 domains with descriptions available. Specialized GO annotations from two levels of
484 hierarchy downstream were taken from files `Domain2GO-Hie-Dist1.csv` and
485 `Domain2GO-Hie-Dist2.csv`. Structural domains functional ontology (SDFO) that
486 mapped information from a theoretic analysis of Domain2GO annotation profiles were
487 reported from the file `SDFO.txt`.

488 *Network construction, visualization and analysis*

489 Mathematical definitions for construction of networks can be found in [Supplementary](#)
490 [Materials and Methods](#). The social network analysis tool Pajek⁴⁸ and the statistical test
491 bench R’s *igraph* package⁴⁹ were used to visualize and analyze the networks, respectively.
492 The collective impact of events was made explicit by Pajek’s Visualization of Similarity
493 (VOS) clustering method^{50,51}. VOS helped reveal communities and design layouts of
494 networks with nodes separated into network modules, where high *modularity indices*
495 ranged from 94-95%. Number of clusters varied over networks (CX, 691; PX, 3,886; PAX,
496 4,126; SPX, 607; SPAX, 620). Network clusters were visually compacted to hubs and their
497 cohesive neighbors with the energy-optimizing Kamada-Kawai ‘separate components’

498 algorithm⁵². Pajek allowed to proportionally reduce the size of highly connected nodes by
 499 some scaling factor for optimally uncluttered visualization. Waterfall and radial network
 500 layouts were designed with node-size scaled down by factors of 0.1 and 0.25, respectively.
 501 R packages equipped with specialized code constructs to draw graphs and derive statistics
 502 were used to analyze network properties^{53,54}. We also used Pre-Hypertext Processing
 503 language (PHP) to write custom scripts that generated radial visualizations of the
 504 networks and helped conduct housekeeping data management⁵⁵. The PHP scripts were
 505 executed in the command line. Results of these scripts were input into Pajek's and R's
 506 analytical procedures. We used the open-source software ImageMagick
 507 (www.imagemagick.org) for batch conversion, captioning, and appending of network
 508 images (to represent legends and scales). A detailed description of partition and data files,
 509 list of network data analysis functions, charting and graphing procedures, methods to
 510 generate power law statistics, modularity indices and randomness checks, and the method
 511 pipeline used to achieve waterfall diagrams can be found in [Supplementary Materials and](#)
 512 [Methods](#).

513 ***Statistical analysis***

514 *Scale-free network behavior*

515 Linear regression models of $P(k)$ given k (i.e. the probability of having k -neighbors) were
 516 used to derive the γ coefficient of the power law distribution and the determination
 517 coefficient, R^2 . The value of γ represents an absolute slope of the log linear model of $P(k)$
 518 vs. k . The slope is usually ≤ 0 . $\gamma \gg 1$ indicates strong tendency towards preferential
 519 attachment. R^2 indicates the percentage of data that fits the linear model. High values of
 520 both γ and R^2 suggest strong scale-free behavior. Additional power law statistics were
 521 calculated as: (i) the exponent of the fitted power law distribution, α , with an assumption
 522 that $P(X=x)$ is proportional to $x^{-\alpha}$; (ii) KS fit statistic to compare the input degree
 523 distribution with that of fitted power-law; and (iii) the KS p -value of a statistical test, with
 524 the null hypothesis that data is being drawn from a power law distribution^{56,57}. $\alpha \gg 1$, $0 <$
 525 $\text{KS fit scores} \ll 1$, and $\text{KS } p\text{-values} \geq 0.05$ suggest that degree data was derived from a
 526 fitted power law distribution. Maximum log likelihood of the fitted scale-free parameters
 527 was also determined. Control networks were included for reference that were generated
 528 with 'Barabási' methods³⁰ of the *igraph* package from R⁴⁹. These controls simulated basic
 529 and extended age-dependent power law graph models given varying sizes of the evolving
 530 networks.

531 *Network modularity*

532 We investigated modularity using six indices: (i) The VOS Quality index (VQ) was
 533 determined using the Pajek VOS algorithm by considering the number or weights of the
 534 links (arcs) between the nodes as similarities. Clusters or communities that were deemed
 535 'similar' were iteratively drawn closer to each other until a final layout was achieved with
 536 least crossings and closest clusters. The quality index VQ was thus calculated for this final
 537 layout as $\sum_{i=1}^c \sum_{j=i+1}^c (e_{ij} - a_i^2)$, where c is the number of communities; e_{ij} is the fraction of

538 edges with one node v in the community i (c_i) and the other node w in the community j
 539 (c_j), defined as $\sum_{vw} (A_{vw}/2m)$ where $v \in c_i$, $w \in c_j$, m is the sum of weights in the graph and
 540 A_{vw} is the weighted value or 0, indicating presence or absence of edge between nodes v
 541 and w in the adjacency matrix A of the network; and a_i is the fraction of weighted k
 542 neighbors attached to the nodes in community i , i.e. $k_i/2m$ ^{50,51}. (ii) The Clustering Ratio
 543 (*C-ratio*) is the ratio of the number of network clusters to the count of the connected
 544 nodes in the network. (iii) The average Clustering Coefficient (C) is defined as the ratio of
 545 the triangles impinging on a node to the connected triples, determined as a global average
 546 over all nodes in a simplified (undirected/unweighted) network^{32,41,42}. C is not meaningful
 547 for strictly bipartite or scale-free graphs⁴⁵. We also report coefficients of linear regression
 548 of C over the age and size of the networks of domain organization. (iv) The Fast Greedy
 549 Community (*FGC*) agglomerative hierarchical algorithm detects community structure for
 550 networks with m edges, n nodes, and a depth d of the dendrogram describing the
 551 community structure, given an optimized linear running time of $O(m \times d \times \log n) \sim$
 552 $O(n \times \log^2 n)$ ⁵⁸. An equivalent modularity index was also calculated using the *Walk Trap*
 553 *Community* (*WTC*) detection algorithm⁵⁹ (results not reported). The *WTC* computation
 554 resembles *FGC* except that *WTC* generates communities using random walks. The
 555 Newman-Girvan algorithm index (NG) was computed with two different input partitions,
 556 the first (v) defined by age (NG_{age}) and the second (vi) defined by VOS clustering (NG_{vos}).
 557 NG calculates the modularity of a network given a predefined division or partition to
 558 measure the influence of the partition in separating the different node types. This
 559 indicates either assortative (positive) or disassortative (negative) mixing across modules³⁷.
 560 The NG algorithm computes an index as $1/(2m) \sum_{ij} (A_{ij} - 1/(2m) k_i k_j) \times \Delta(c_i, c_j)$, where m is the
 561 sum total of weights in the graph and A_{ij} are weighted entries in the adjacency matrix of
 562 the network; $k_i \mid k_j$ and $c_i \mid c_j$ are the weighted degrees and the components (numeric
 563 partitions) of the nodes i and j , respectively; finally, $\Delta(x, y)$ equals 1 if $x=y$ and 0
 564 otherwise³⁷. We also computed the NG index for two additional input memberships
 565 generated by *FGC* and *WTC* (results not reported). The VQ , *C-ratio*, C and *FGC* indices
 566 each range from 0 to 1, while the NG indices range from -1 to 1. In all cases, higher values
 567 represent strong modularity of the network at an event of evolutionary history. Heatmaps
 568 of modularity were constructed using log10-scaled modularity matrices, with each map
 569 element given as $(A_{ij} - k_i k_j / (2m)) M_{nd}$, where A_{ij} , k_i , k_j and m were the same as defined for
 570 NG ³⁷, while M_{nd} was the network's modularity index at event nd . Cladistic representations
 571 of modularity were visualized with dendrograms whose metrics were calculated from
 572 squared Euclidean distance matrices, which indicate dissimilarities between cluster
 573 means⁶⁰. The dissimilarity or distance matrices were clustered hierarchically using the
 574 Ward's minimum variance method that seeks compact and spherical clusters⁶¹.

575 *Quantifying randomness in networks*

576 The Bartels rank test of randomness, which primarily offers a rank version of von
 577 Neumann's Ratio Test for Randomness⁶², was used to measure random network behavior.
 578 The resultant test statistic RVN is defined as $\sum_{i=1 \rightarrow n-1} (R_i - R_{i+1})^2 / \sum_{i=1 \rightarrow n} (R_i - (n+1)/2)^2$,
 579 where $R_i = \text{rank}(X_i)$ with $i=1 \dots n$, $(RVN-2)/\sigma$ is the asymptotically standard normal,
 580 and $\sigma^2 = [4(n-2)(5n^2-2n-9)]/[5n(n+1)(n-1)^2]$. The null hypothesis of this method was

581 randomness, which was tested against the alternate hypothesis of non-randomness, given
 582 a trend of RVN values. A p-value is computed from a two-sided beta distribution
 583 approximation test. Random graph controls were created by following the Erdős–Rényi
 584 graph model^{33,63}.

585 References

- 586 1. Chothia, C. & Gough, J. Genomic and structural aspects of protein evolution.
 587 *Biochem. J* **419**, 15–28 (2009).
- 588 2. Chothia, C., Gough, J., Vogel, C. & Teichmann, S. A. Evolution of the protein
 589 repertoire. *Science (80-.)*. **300**, 1701 (2003).
- 590 3. Nath, N., Mitchell, J. B. O. & Caetano-Anollés, G. The natural history of
 591 biocatalytic mechanisms. *PLoS Comput Biol* **10**, e1003642 (2014).
- 592 4. Debès, C., Wang, M., Caetano-Anollés, G. & Gräter, F. Evolutionary optimization
 593 of protein folding. *PLoS Comput Biol* **9**, e1002861 (2013).
- 594 5. Aziz, M. F. *et al.* Stress induces biphasic-rewiring and modularization patterns in
 595 the metabolomic networks of *Escherichia coli*. *IEEE Intl. Conf. Bioinf. Biomed.* 593–
 596 597 (2012) doi:10.1109/BIBM.2012.6392626.
- 597 6. Aziz, M. F., Caetano-Anollés, K. & Caetano-Anollés, G. The early history and
 598 emergence of molecular functions and modular scale-free network behavior. *Sci.*
 599 *Rep.* **6**, (2016).
- 600 7. Apic, G., Gough, J. & Teichmann, S. A. Domain combinations in archaeal,
 601 eubacterial and eukaryotic proteomes. *J. Mol. Biol.* **310**, 311–325 (2001).
- 602 8. Wang, M. & Caetano-Anollés, G. The evolutionary mechanics of domain
 603 organization in proteomes and the rise of modularity in the protein world.
 604 *Structure* **17**, 66–78 (2009).
- 605 9. Murzin, A. G., Brenner, S. E., Hubbard, T., Chothia, C. & others. SCOP: a
 606 structural classification of proteins database for the investigation of sequences and
 607 structures. *J. Mol. Biol.* **247**, 536–540 (1995).
- 608 10. Fox, N. K., Brenner, S. E. & Chandonia, J.-M. SCOPe: Structural Classification of
 609 Proteins—extended, integrating SCOP and ASTRAL data and classification of new
 610 structures. *Nucleic Acids Res.* **42**, D304–D309 (2013).
- 611 11. Gerstein, M. How representative are the known structures of the proteins in a
 612 complete genome? A comprehensive structural census. *Fold. Des.* **3**, 497–512
 613 (1998).
- 614 12. Wang, M., Kurland, C. G. & Caetano-Anollés, G. Reductive evolution of
 615 proteomes and protein structures. *Proc. Natl. Acad. Sci.* **108**, 11954–11958 (2011).
- 616 13. Illergård, K., Ardell, D. H. & Elofsson, A. Structure is three to ten times more
 617 conserved than sequence—a study of structural response in protein cores. *Proteins*
 618 *Struct. Funct. Bioinforma.* **77**, 499–508 (2009).

- 619 14. Bashton, M. & Chothia, C. The geometry of domain combination in proteins1. *J.*
620 *Mol. Biol.* **315**, 927–939 (2002).
- 621 15. Vogel, C., Berzuini, C., Bashton, M., Gough, J. & Teichmann, S. A. Supra-domains:
622 evolutionary units larger than single protein domains. *J. Mol. Biol.* **336**, 809–823
623 (2004).
- 624 16. Caetano-Anollés, G. *et al.* The origin and evolution of modern metabolism. *Int. J.*
625 *Biochem. Cell Biol.* **41**, 285–297 (2009).
- 626 17. Caetano-Anollés, G. & Caetano-Anollés, D. An evolutionarily structured universe
627 of protein architecture. *Genome Res.* **13**, 1563–1571 (2003).
- 628 18. Wang, M. & Caetano-Anollés, G. Global phylogeny determined by the
629 combination of protein domains in proteomes. *Mol. Biol. Evol.* **23**, 2444–2454
630 (2006).
- 631 19. Wang, M., Yafremava, L. S., Caetano-Anollés, D., Mittenthal, J. E. & Caetano-
632 Anollés, G. Reductive evolution of architectural repertoires in proteomes and the
633 birth of the tripartite world. *Genome Res.* **17**, 1572–1585 (2007).
- 634 20. Wang, M. *et al.* A universal molecular clock of protein folds and its power in
635 tracing the early history of aerobic metabolism and planet oxygenation. *Mol. Biol.*
636 *Evol.* **28**, 567–582 (2011).
- 637 21. Caetano-Anollés, D., Kim, K. M., Mittenthal, J. E. & Caetano-Anollés, G. Proteome
638 evolution and the metabolic origins of translation and cellular life. *J. Mol. Evol.* **72**,
639 14–33 (2011).
- 640 22. Mittenthal, J. E., Caetano-Anollés, D. & Caetano-Anollés, G. Biphasic patterns of
641 diversification and the emergence of modules. *Front. Genet.* **3**, 147 (2012).
- 642 23. MacDougall, M. H. Simulating computer systems: Techniques and tools. (1987).
- 643 24. Delaney, W. & Vaccari, E. Dynamic models and discrete event simulation. (1989).
- 644 25. Pidd, M. Computer simulation in management science. *JOURNAL-*
645 *OPERATIONAL Res. Soc.* **57**, 327 (2006).
- 646 26. Caetano-Anollés, G., Kim, H. S. & Mittenthal, J. E. The origin of modern metabolic
647 networks inferred from phylogenomic analysis of protein architecture. *Proc. Natl.*
648 *Acad. Sci.* **104**, 9358 (2007).
- 649 27. Caetano-Anollés, G., Kim, K. M. & Caetano-Anollés, D. The phylogenomic roots
650 of modern biochemistry: Origins of proteins, cofactors and protein biosynthesis. *J.*
651 *Mol. Evol.* **74**, 1–34 (2012).
- 652 28. Caetano-Anollés, K. & Caetano-Anollés, G. Structural phylogenomics reveals
653 gradual evolutionary replacement of abiotic chemistries by protein enzymes in
654 purine metabolism. *PLoS One* **8**, e59300 (2013).
- 655 29. Apic, G., Huber, W. & Teichmann, S. A. Multi-domain protein families and
656 domain pairs: comparison with known structures and a random model of domain

- 657 recombination. *J. Struct. Funct. Genomics* **4**, 67–78 (2003).
- 658 30. Barabási, A. L. & Albert, R. Emergence of scaling in random networks. *Science* (80-.
659). **286**, 509 (1999).
- 660 31. Pang, T. Y. & Maslov, S. Universal distribution of component frequencies in
661 biological and technological systems. *Proc. Natl. Acad. Sci.* **110**, 6235–6239 (2013).
- 662 32. Ravasz, E., Somera, A. L., Mongru, D. A., Oltvai, Z. N. & Barabási, A. L.
663 Hierarchical organization of modularity in metabolic networks. *Science* (80-.). **297**,
664 1551–1555 (2002).
- 665 33. Erdős, P. & Rényi, A. Connectivity of random nets. *Publ. Math. Inst. Hungarian*
666 *Acad. Sci.* **5**, 17–61 (1960).
- 667 34. Bollobas, B. *Random Graphs*. (Academic Press, London, 1985).
- 668 35. Strogatz, S. H. Exploring complex networks. *Nature* **410**, 268–276 (2001).
- 669 36. Mughal, F. & Caetano-Anollés, G. MANET 3.0: Hierarchy and modularity in
670 evolving metabolic networks. *PLoS One* **14**, e0224201 (2019).
- 671 37. Newman, M. E. J. & Girvan, M. Finding and evaluating community structure in
672 networks. *Phys. Rev. E* **69**, 26113 (2004).
- 673 38. Jeong, H., Tombor, B., Albert, R., Oltvai, Z. N. & Barabási, A. L. The large-scale
674 organization of metabolic networks. *Nature* **407**, 651–654 (2000).
- 675 39. Overbeek, R. *et al.* WIT: integrated system for high-throughput genome sequence
676 analysis and metabolic reconstruction. *Nucleic Acids Res.* **28**, 123–125 (2000).
- 677 40. Wagner, A. & Fell, D. A. The small world inside large metabolic networks. *Proc. R.*
678 *Soc. London. Ser. B Biol. Sci.* **268**, 1803–1810 (2001).
- 679 41. Wasserman, S. & Faust, K. *Social network analysis: Methods and applications.* **8**,
680 (1994).
- 681 42. Barrat, A., Barthelemy, M., Pastor-Satorras, R. & Vespignani, A. The architecture
682 of complex weighted networks. *Proc. Natl. Acad. Sci. U. S. A.* **101**, 3747–3752
683 (2004).
- 684 43. Watts, D. J. & Strogatz, S. H. Collective dynamics of ‘small-world’ networks. *Nature*
685 **393**, 440–442 (1998).
- 686 44. Albert, R. & Barabási, A. L. Statistical mechanics of complex networks. *Rev. Mod.*
687 *Phys.* **74**, 47 (2002).
- 688 45. Newman, M. E. J., Strogatz, S. H. & Watts, D. J. Random graphs with arbitrary
689 degree distributions and their applications. *Phys. Rev. E* **64**, 26118 (2001).
- 690 46. Caetano-Anollés, G., Wang, M. & Caetano-Anollés, D. Structural phylogenomics
691 retrodicts the origin of the genetic code and uncovers the evolutionary impact of
692 protein flexibility. *PLoS One* **8**, e72225 (2013).
- 693 47. Tal, G., Boca, S. M., Mittenthal, J. & Caetano-Anollés, G. A dynamic model for the

- 694 evolution of protein structure. *J. Mol. Evol.* **82**, 230–243 (2016).
- 695 48. Mrvar, A. & Batagelj, V. Analysis and visualization of large networks with program
696 package Pajek. *Complex Adapt. Syst. Model.* **4**, 1–8 (2016).
- 697 49. Csardi, G. & Nepusz, T. The igraph software package for complex network
698 research. *InterJournal, Complex Syst.* **1695**, 1–9 (2006).
- 699 50. Van Eck, N. J. & Waltman, L. VOS: a new method for visualizing similarities
700 between objects. in *Advances in Data Analysis: Proceedings of the 30th Annual*
701 *Conference of the German Classification Society* 299–306 (Heidelberg: Springer
702 Verlag, 2007).
- 703 51. Waltman, L., van Eck, N. J. & Noyons, E. C. M. A unified approach to mapping
704 and clustering of bibliometric networks. *J. Informetr.* **4**, 629–635 (2010).
- 705 52. Kamada, T. & Kawai, S. An algorithm for drawing general undirected graphs. *Inf.*
706 *Process. Lett.* **31**, 7–15 (1989).
- 707 53. Ihaka, R. & Gentleman, R. R: A language for data analysis and graphics. *J. Comput.*
708 *Graph. Stat.* **5**, 299–314 (1996).
- 709 54. R Core Team. R: A language and environment for statistical computing. R
710 Foundation for Statistical Computing, Vienna, Austria. 2013. (2014).
- 711 55. PHP-Group & others. PHP: Hypertext PreProcessor. *Internet <http://www.php.net>*
712 (2012).
- 713 56. Newman, M. E. J. Power laws, Pareto distributions and Zipf's law. *Contemp. Phys.*
714 **46**, 323–351 (2005).
- 715 57. Clauset, A., Shalizi, C. R. & Newman, M. E. J. Power-law distributions in empirical
716 data. *SIAM Rev.* **51**, 661–703 (2009).
- 717 58. Clauset, A., Newman, M. E. J. & Moore, C. Finding community structure in very
718 large networks. *Phys. Rev. E* **70**, 66111 (2004).
- 719 59. Pons, P. & Latapy, M. Computing communities in large networks using random
720 walks. *J. Graph Algorithms Appl.* **10**, 191–218 (2006).
- 721 60. Borg, I. & Groenen, P. Modern multidimensional scaling: theory and applications.
722 *J. Educ. Meas.* **40**, 277–280 (2003).
- 723 61. Murtagh, F. & Legendre, P. Ward's hierarchical agglomerative clustering method:
724 Which algorithms implement ward's criterion? *J. Classif.* **31**, 274–295 (2014).
- 725 62. Bartels, R. The rank version of von Neumann's ratio test for randomness. *J. Am.*
726 *Stat. Assoc.* **77**, 40–46 (1982).
- 727 63. Erdős, P. & Rényi, A. On random graphs, I. *Publ. Math.* **6**, 290–297 (1959).

728 **Acknowledgments**

729

730 Research was supported by grants from the National Science Foundation (MCB-0749836
731 and OISE-1132791) and the United States Department of Agriculture (ILLU-802-909 and
732 ILLU-483-625) to GCA. Materials and data necessary to interpret the findings of this
733 paper have been included in the manuscript.

734 **Author Contributions**

735

736 G.C.-A. conceptualized the study. M.F.A. generated primary data, conducted network
737 analysis, and generated figures and written documentation. Both authors interpreted
738 results and wrote and revised the manuscript.

739 **Additional Information**

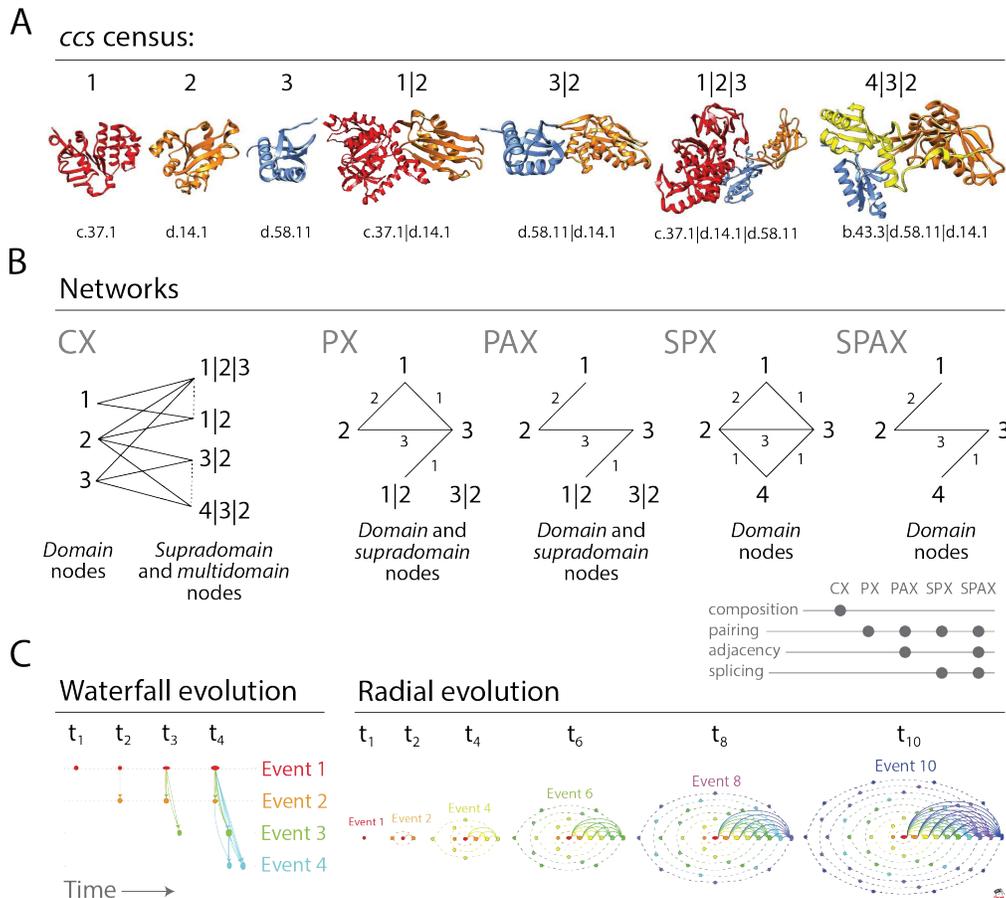
740

741 **Supplementary information** accompanies this paper at <http://www.nature.com/srep>.

742

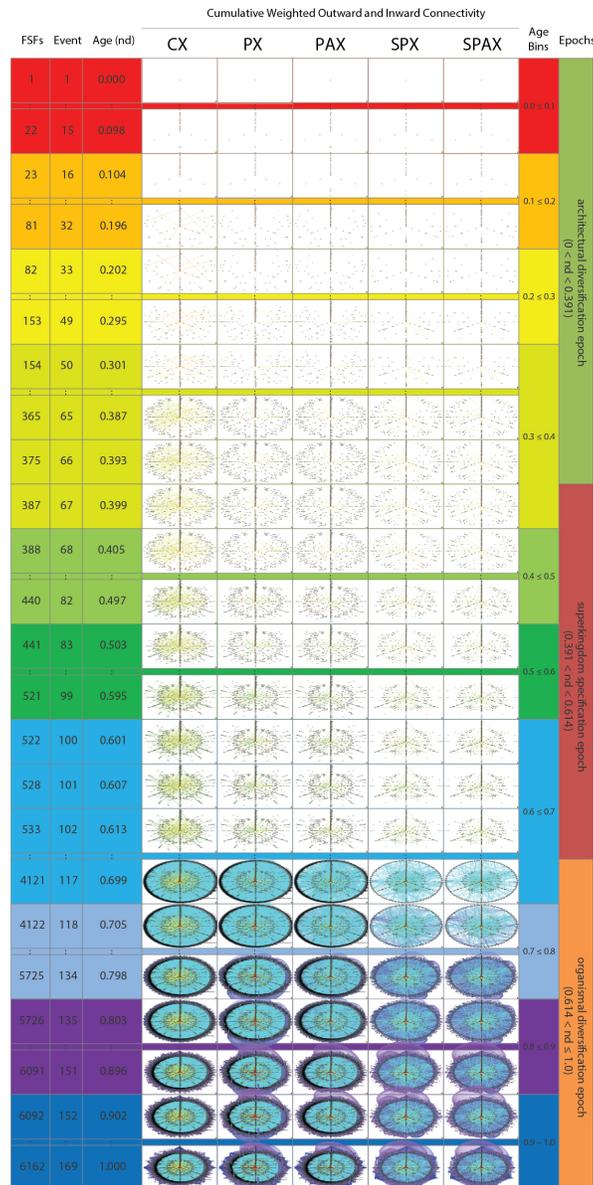
743 **Competing financial interests:** The authors declare no competing financial interests.

744



745

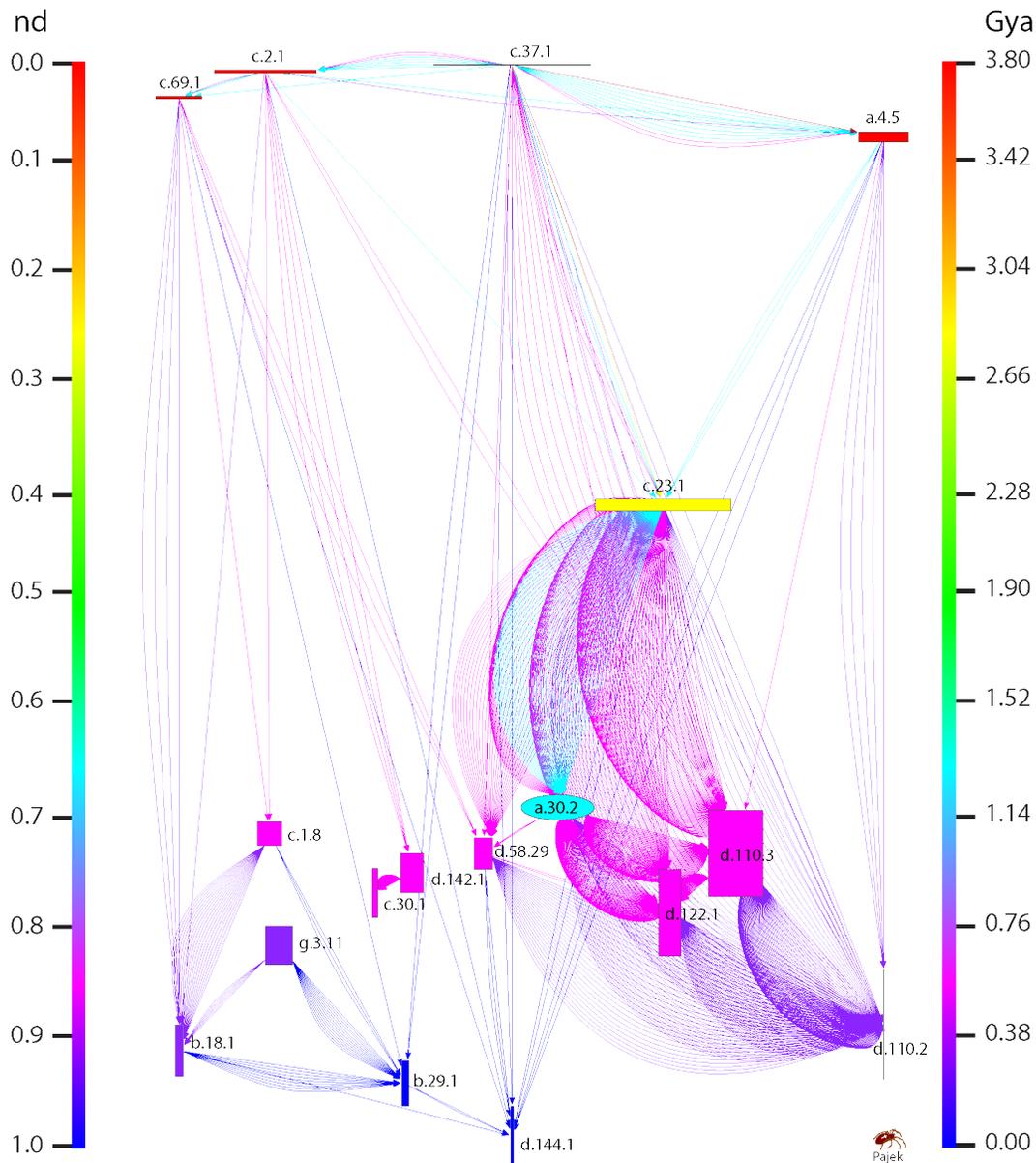
746 **Figure 1. Networks of protein domain organization.** (A) The genomic census of structural domains and their
 747 combinations defines SCOP *concise classification string* (*ccs*) descriptors of domains, supradomains and multidomains
 748 that are building blocks of networks. We illustrate the census with a sample from the entire entity set, comprising of 3
 749 domains (1, 2 and 3), 2 supradomains (1|2 and 3|2) and 2 multidomains (1|2|3 and 4|3|2) that are common in
 750 dehydratase enzymes and elongation factors. *ccs* identifiers of structural domain constituents defined at fold
 751 superfamily (FSF) level are listed below the atomic models visualized in ribbon format with Chimera. (B) Five operative
 752 criteria for network generation capture the interactions among protein architecture nodes as networks grow in
 753 evolution. CX is a partial bipartite network (projection-decomposable) that connects domain nodes to supradomain
 754 and multidomain nodes (which can connect to each other; hatched links) when present in multidomain proteins. PX
 755 connects domain and supradomain nodes when multidomain proteins are ‘decomposed’ into pairs of architectures,
 756 regardless of topological constraints. PAX borrows the PX criterion but respects topological constraints. SPX connects
 757 domain nodes spliced from architectures when domain pairs are present in proteins. SPAX connects domain nodes
 758 when adjacent domain pairs are present in proteins. (C) Chronological development of evolving networks. In ‘waterfall
 759 evolution’ layout, time progresses from left to right as ‘discrete events’ of network evolution progressively unfold the
 760 appearance of nodes and links (time-directed arrows known as arcs) from top to bottom, colored according to their age.
 761 Arc multiplicities describe link cardinality. Source-sink recruitments of architectures are visualized by horizontal and
 762 vertical elongations of node symbols, which describe their outdegree and indegree, respectively. As networks grow, the
 763 symbols of older nodes widen by outdegree accumulation, while those of younger nodes grow tall by indegree
 764 accumulation. In ‘radial evolution’ layout, the time-variant network grows by accumulating nodes in concentric rings
 765 (orbitals), each reflecting a time event. We illustrate radial evolution with 6 snapshots of a network growing to a size of
 766 55 nodes as it unfolds from time t₁ to t₁₀. Nodes (n) in orbitals (r) grow at r+1 rate and only one node per orbital
 767 connects to single nodes in each of the other orbitals. Thus, outward links (o) of an orbital are o=t-r-1, where t is the
 768 current time. Inward links (i) of an orbital are i=t-o-1=r. Finally, total links of a network at any time are t(t-1)/2. The
 769 width and height of symbols represent the outdegree and indegree of nodes, respectively. Symbol sizes are shifted by 10
 770 for a better visualization of nodes.



771

772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787

Figure 2. Evolving networks in radial evolution layout. Snapshots of network growth describe the evolution of 6,162 domain, supradomain and multidomain architectures or 1,643 domains spliced from them. They represent 24 out of 169 time events of the evolutionary timeline, which are indexed with evolutionary age (nd, ranging from 0.0 to 1.0), age bin (one of 10), and one of the 3 epochs of protein evolution (Wang *et al.*, 2007). Age bins were custom RGB color-coded to highlight the flow of time, from top to bottom. The evolving CX, PX, PAX, SPX and SPAX networks reveal the gradual evolutionary accumulation of nodes and links. The sizes of the horizontal and vertical axes of the node symbols depict outward and inward weighted connectivity, respectively, with all weighted degree vectors shifted by 10 for visualization and inclusion of 0-degree nodes. The curved arcs describe recurring interactions between architectures that are accumulating along the successive events of the timeline. Arcs symbolize the flow of time from ancient to recent architectures and are color-coded according to the age of the more recent of the component nodes involved; arcs between contemporary nodes are excluded. Since, in pairwise networks the age of the most recent parent node could be assigned to the arc, the connectivity-defining pairing events are absent in the first (red) and the first and second (red, orange) bins of the PX and SPX and the PAX and SPAX networks, respectively. The angles of multiple arcs emerging from nodes are incremented by 2 to avoid overlap. Node RGB colors represent age. Grey-scale color of node borders depict fusional/fissional properties (Supplementary Fig. S3). Node shapes describe GO categories: circle, molecular function; squares, biological process; rhomboid, cellular component; triangle, unassigned.



788

789

790

791

792

793

794

795

796

797

798

799

800

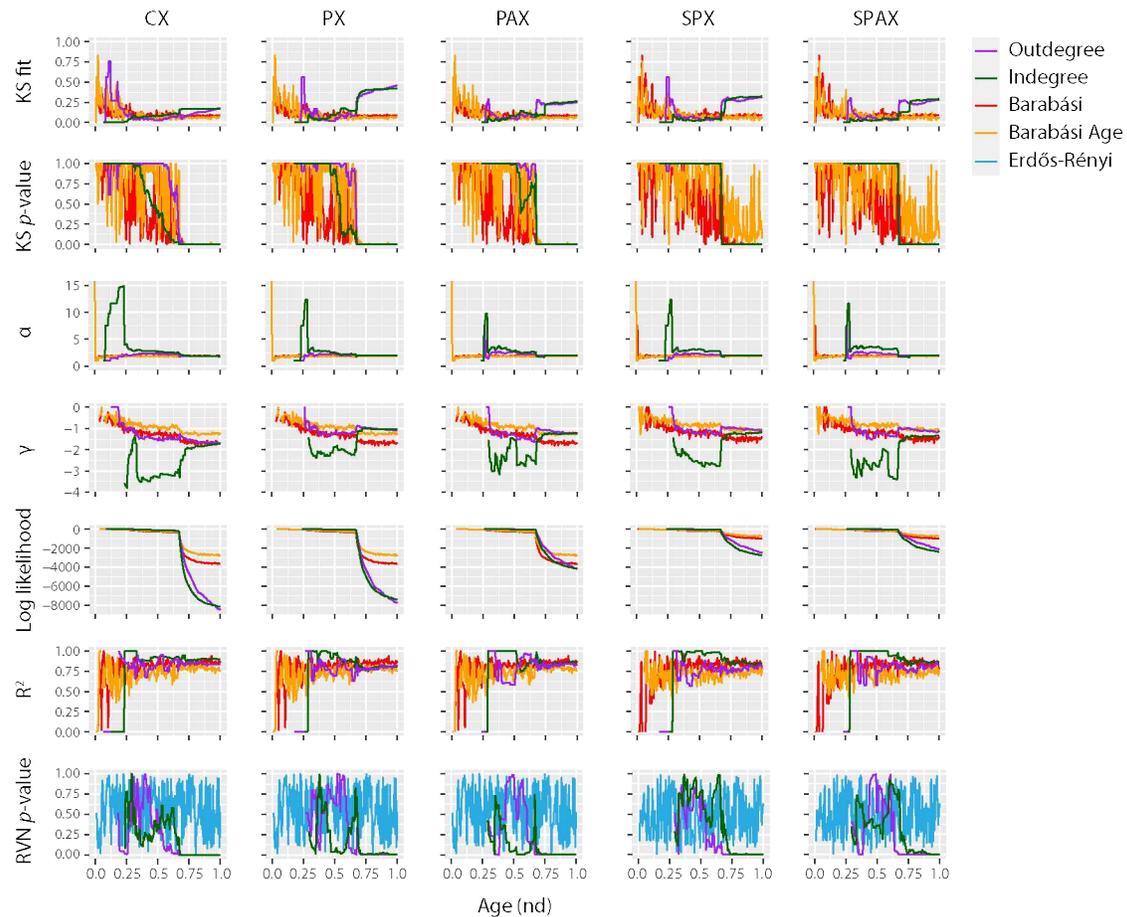
801

802

803

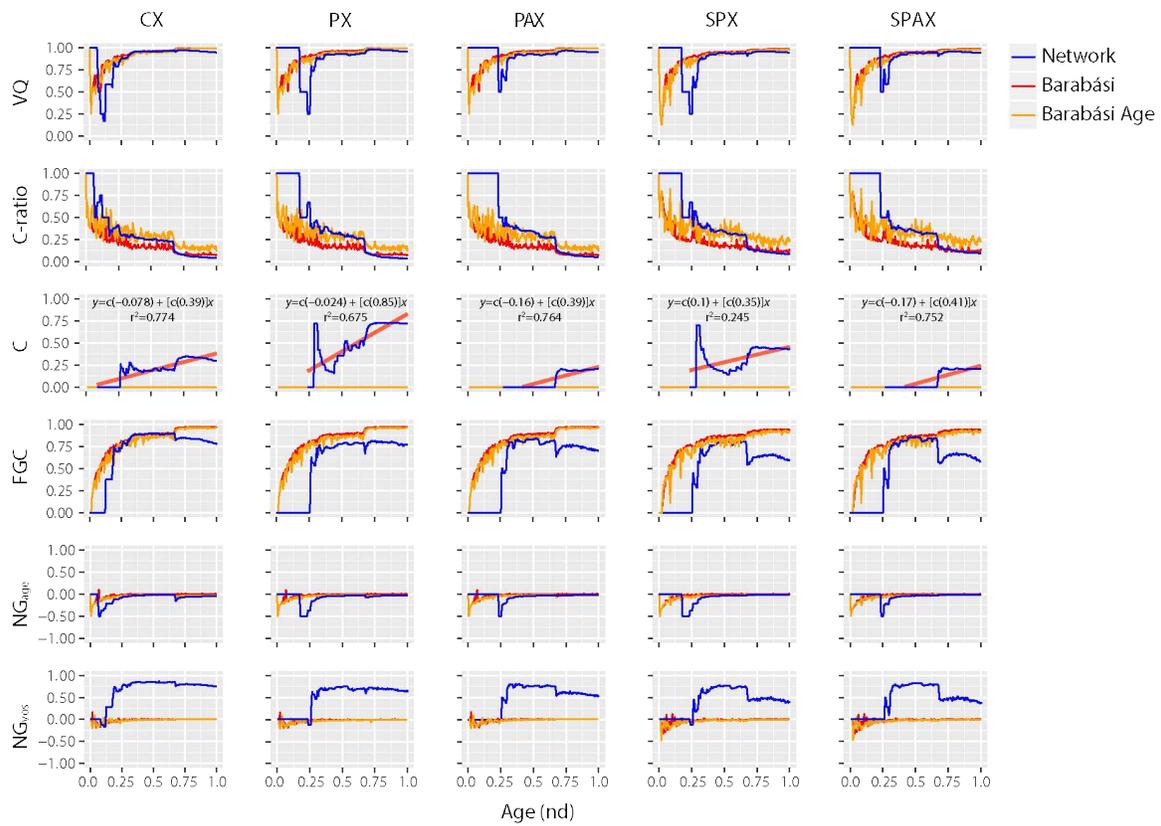
Figure 3. An extant SPX network in waterfall layout describing the evolution of spliced domains with the largest (100th percentile) network connectivity. The SPX network of 1,643 spliced domains was reduced with the restrictive criterion of excluding nodes with combined outdegrees and indegrees $\geq 99\%$ of those of the rest of the nodes. The set of arcs (arched arrows symbolizing flow of time) was also reduced to pairing events between domains in the 100th percentile connectivity and excluded those between contemporary nodes. Nodes are arranged top-down and colored according to age (*nd*) on a relative 0-to-1 scale that describes evolutionary time events. Ages are also time-calibrated with a molecular clock of FSF domains, which uses fossils and microfossils, geochemical, biochemical, and biomarker data²⁰. FSF origin is given in billion years ago (Gya). Nodes were labeled with SCOP ccs domain descriptors. To showcase source-and-sink relationships, node symbol sizes were scaled proportional to the weighted outdegree and indegree along the horizontal and vertical axes, respectively. Weighted degrees were scaled as $\times 2+2$ to include 0-degree nodes for better visualization. The modular spread of nodes was based on VOS clustering (see methods). Arcs are color-coded according to the age of the more recent of the component nodes involved; no arcs were present in the ancient-most age bin (red) of the timeline. Angles of multiple arcs emerging from nodes are incremented by 2 to avoid overlap. See caption of Figure 2 for indexing of node colors and shapes.

804



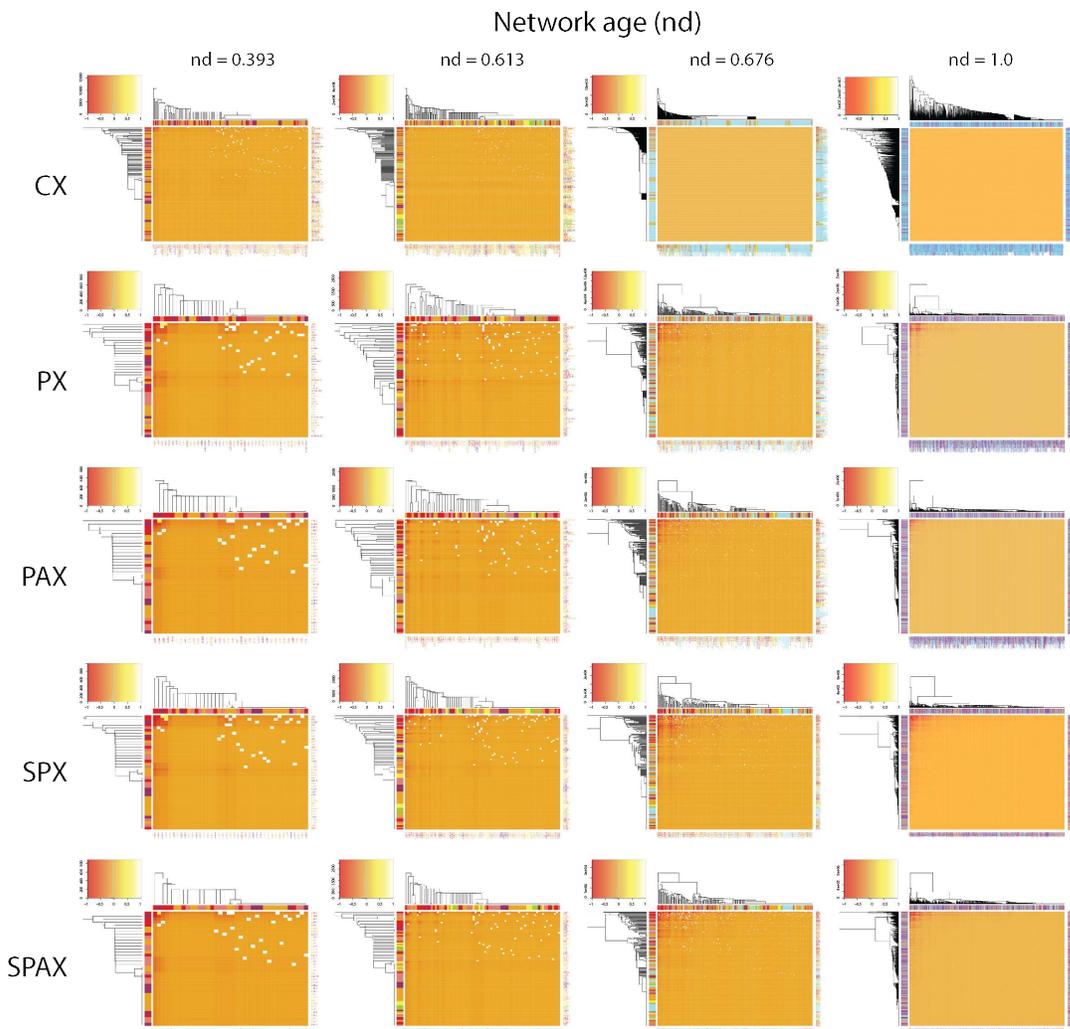
805

806 **Figure 4. Statistical descriptors of power law and random behavior.** Six indicators of preferential attachment were
 807 studied along the evolutionary timeline to explore processes of network growth, with network age (nd) indicated on a
 808 relative 0-to-1 scale. Outdegree and indegree connections were cumulative and weighted in evolving networks. Barabási
 809 (red) and Barabási-Age (orange) networks were included as control sets. The Barabási model specifies the probability of
 810 preference of an old node as $P_i \sim k_i^\alpha$ while the Barabási-Age model grants heavier power law properties to older nodes
 811 (exhibiting smaller nd) with $P_i \sim (k_i^\alpha)(l_i^\beta)$, where k_i is the indegree of node i of the current event, α is the preferential
 812 attachment exponent ($\alpha = 1$ for linear preferential attachment), l_i is the age of node i , i.e. the number of events elapsed
 813 since the node was added, with maximum number measured by the 'aging.bin' parameter, and β is the aging exponent
 814 ($\beta = 1$ for linear increases in probability of preference of an older node with high l_i). Power law indices include: (i) the
 815 KS fit statistic that compares the input degree data distribution with the fitted power law distribution (smaller scores
 816 denote better fit); (ii) the KS p-value, which rejects the null hypothesis that degree data was drawn from the fitted
 817 power-law distribution when less than $\alpha=0.05$; (iii) the exponent of the fitted power-law distribution (α); (iv) the slope
 818 of power-law linear regression model (γ); (v) the log-likelihood of the fitted parameters; and (vi) the coefficient of
 819 determination (R^2) that measures the percentage of degree data that fits the linear model. The randomness of the
 820 evolving networks was quantified by the p-value of an approximated beta distribution from the rank version of von
 821 Neumann's Ratio Test for Randomness⁶² ($RVN_{p\text{-value}}$). The alternate hypothesis was non-randomness. Comparative
 822 graphs of strictly random Erdős-Rényi control networks of corresponding sizes at the given time-events were also
 823 plotted. Lower KS fit, higher KS p-value, higher α , lower γ and near-zero likelihood, given lower $RVN_{p\text{-value}}$, support
 824 power law behavior.



825

826 **Figure 5. Network modularity.** Six indicators of modularity were studied along the evolutionary timeline to explore
 827 the evolution of network structure, with network age (*nd*) indicated on a relative 0-to-1 scale. Modularity indices
 828 include the VOS Quality (VQ) index, the Clustering ratio (C-ratio), the average Clustering Coefficient (C), the Fast
 829 Greedy Community (FGC) index, and the Newman-Girvan index defined by age (NG_{age}) or VOS clustering (NG_{vos}).
 830 Modularity calculations required cumulative, undirected, and weighted connectivity input. The Barabási (red) and
 831 Barabási-Age (orange) models (see caption of Figure 4) were included as control sets. The regressions of C with age
 832 (*nd*) are shown as linear models (red lines) for each network together with supporting determination coefficients (R²).



833

834 **Figure 6. Evolution of modularity and hierarchical organization of networks over select events of the evolutionary**
835 **timeline.** NG_{age} pairwise modularity values³⁹, scaled by \log_{10} of network-wide absolute modularity values, were used as
836 input for the calculation of Euclidean distance matrices⁵⁴, which were visualized as heatmaps. Heatmap tiles represent
837 modular strength between any two architectures relative to the respective strength of their linkages to other
838 architectures of the network. The embedded dendrograms that define the order of rows and columns of the heatmaps
839 were generated by hierarchical clustering of the distance matrices with the Ward's minimum variance method⁵⁵. The
840 height of dendrograms represents dissimilarity between clusters while the clades show grouping rearrangements of
841 architectures. The top-left insets depict frequency histograms of the heatmap modularity values scaled from -1 to 1 (i.e.
842 disassortative to assortative). The four panels describe growth of each evolving network (left-to-right). Network age
843 corresponds to the middle approximate boundaries of the three evolutionary epochs of the protein world
844 (Supplementary Fig. S2), i.e., end of 'architectural diversification' ($nd = 0.393$), end of 'superkingdom specification' (nd
845 $= 0.613$), onset of the 'big bang' of domain organization at the start of 'organismal diversification' ($nd = 0.676$); and the
846 present ($nd = 1$). Nodes were age-sorted ascendingly within clusters and labelled using standard SCOP nomenclature¹⁷.
847 In the case of SPX and SPAX, nodes correspond to 1,643 domains mapped to the entity set of 6,162 architectures. The
848 color-coding of bands and labels identifies the age of architectures (Supplementary Fig. S2). The relatively 'flatter'
849 heatmap and 'skewed' dendrogram patterns of CX (typically at $nd = 0.667$ and $nd = 1.000$) are an artifact of unweighted
850 distance matrices of CX, which contrast with the weighted ones of pairwise criterion-based networks. The most
851 prominent clades correspond to the modules of the most ancient domain structures harboring the two major waves of
852 architectural innovation. We also generated heatmaps of power-law control networks of corresponding sizes at the
853 given time-events (Supplementary Fig. S10). When compared to the pairwise networks, the combined heatmap and
854 dendrogram patterns of CX suggest a hidden switch from scale-freeness to modular behavior, eventually giving rise to
855 hierarchical modularity with visible emergence of modules within modules.

856
857
858

Table 1: Domains and domain combinations scoring $\geq 99.9^{\text{th}}$ percentiles of 249,916, [64] and {23}, based on combined outdegrees of the five networks at time points 1.0, [0.676] and {0.671}, respectively. The square and curly brackets denote values from the events after and before the big bang, respectively.

Age Rank	Label	Node Age	Network(s)	Out Degree	Fussional / Fissional	Description	GO Name
388	c.23.1	0.4046243	PX, PAX, SPX	1013, 390, 330	fissional/fussional	CheY-like	regulation of multicellular organismal development
1	c.37.1	0.0000000	PX, SPX, PAX, SPAX, CX	607, 380, 376, 314, 271	fussional	P-loop containing nucleoside triphosphate hydrolases	positive regulation of reproductive process
			[CX, PX, SPX, PAX, SPAX]	[109, 97, 74, 64, 64]			
			{CX, PX, SPX, PAX, SPAX}	{34, 32, 26, 23, 23}			
2446	a.30.2	0.6820809	PX	578	fissional/fussional	Homodimeric domain of signal transducing histidine =	alkene binding
48	e.23.1	0.1445087	PX	452	fussional	Acetyl-CoA synthetase-like	regulation of primary metabolic process
			[PX]	[75]			
2	c.2.1	0.0057803	PX	427	fussional	NAD(P)-binding Rossmann-fold domains	pyridine-containing compound metabolic process
			[PX, CX, SPX]	[117, 80, 66]			
1518	d.110.3 a.30.2	0.6763006	PX	423	fissional/fussional	#N/A	#N/A
283	a.28.1	0.3526012	PX	416	fussional	ACP-like	cell periphery
2543	d.110.3&	0.6820809	PX	369	fissional/fussional	#N/A	#N/A
858	d.110.2 d.110.3	0.6763006	PX	357	fissional/fussional	#N/A	#N/A
187	c.30.1 d.142.1	0.3179191	PX	352	fissional/fussional	#N/A	#N/A
8	c.69.1	0.0346821	PX	327	fussional	alpha/beta-Hydrolases	regulation of multicellular organismal development
4777	d.110.3	0.7225434	PX	325	fissional/fussional	PYP-like sensor domain (PAS domain)	regulation of cellular macromolecule biosynthetic process
17	c.23.16	0.0809249	PX	315	fussional	Class I glutamine amidotransferase-like	positive regulation of oxidative phosphorylation uncoupler activity
4465	d.122.1 c.23.1	0.7109827	PX	291	fussional/fissional/fussional	#N/A	#N/A
1599	d.110.3 d.110.2	0.6763006	PX	262	fissional/fussional	#N/A	#N/A
443	c.1.33	0.5028902	PX	253	fussional	EAL domain-like	cyclic-guanylate-specific phosphodiesterase activity

859
860

861
862
863

Table 2: Domains and domain combinations scoring $\geq 99.9^{\text{th}}$ percentile of 247.977, [21] and {5}, based on combined indegrees of the five networks at time points 1.0, [0.676] and {0.671}, respectively. The square and curly brackets denote values from the events after and before the big bang, respectively.

Age Rank	Label	Node Age	Network(s)	In Degree	Fusional / Fissional	Description	GO Name
6044	d.110.2	0.8728324	PX, PAX, SPX	766, 295, 267	fissional	GAF domain-like	purine-containing compound catabolic process
4777	d.110.3	0.7225434	PX	735	fissional/fusional	PYP-like sensor domain (PAS domain)	regulation of cellular macromolecule biosynthetic process
5529	d.122.1	0.7745665	PX	701	fissional/fusional	ATPase domain of HSP90 chaperone/DNA topoisomerase =	nucleic acid metabolic process
5038	a.30.2 d.122.1	0.7341040	PX	550	fusional/fissional/fusional	#N/A	#N/A
5101	c.43.1	0.7398844	PX	445	fissional/fusional	CoA-dependent acyltransferases	monocarboxylic acid catabolic process
5664	d.110.3 a.30.2 d.122.1	0.7919075	PX	439	fusional/fissional	#N/A	#N/A
6150	c.43.1&	0.9768786	PX	432	fusional/fissional	#N/A	#N/A
5304	c.30.1	0.7572255	PX	375	fissional/fusional	PreATP-grasp domain	pyrimidine-containing compound biosynthetic process
6148	b.1.1	0.9768786	PX	370	fissional	Immunoglobulin	regulation of mesoderm development
4848	e.23.1 a.28.1	0.7225434	PX	367	fusional/fissional	#N/A	#N/A
5095	d.142.1	0.7398844	PX	359	fissional/fusional	Glutathione synthetase ATP-binding domain-like	pyrimidine-containing compound biosynthetic process
5731	g.3.11	0.8034682	PX	317	fissional/fusional	EGF/Laminin	positive regulation of receptor activity
4118	c.43.1& e.23.1 a.28.1	0.6994219	PX	287	fusional/fissional/fusional	#N/A	#N/A
4758	d.58.29	0.7225434	PX	272	fissional/fusional	Nucleotide cyclase	regulation of primary metabolic process
5521	d.110.3 d.58.29	0.7745665	PX	266	fusional/fissional	#N/A	#N/A
4855	c.43.1& e.23.1 a.28.1 c.43.1&	0.7225434	PX	265	fusional/fissional	#N/A	#N/A
4763	a.30.2 d.122.1 c.23.1	0.7225434	PX	263	fusional/fissional/fusional	#N/A	#N/A
5768	b.23.1	0.8092486	PX	261	fissional/fusional	Spermadhesin, CUB domain	regulation of anatomical structure size
5759	b.1.2	0.8092486	PX	259	fissional/fusional	Fibronectin type III	regulation of CD4-positive, alpha-beta T cell activation
2886	c.43.1 e.23.1 a.28.1	0.6878613	PX	258	fusional/fissional/fusional	#N/A	#N/A
[1620]	c.43.1& e.23.1	0.6763006	PX	28	fusional/fissional/fusional	#N/A	#N/A
[1223]	d.142.1 c.24.1	0.6763006	PX	25	fusional/fissional/fusional	#N/A	#N/A
[1311]	a.28.1 c.43.1&	0.6763006	PX	25	fusional/fissional/fusional	#N/A	#N/A
[1032]	e.23.1 a.28.1 c.43.1& e.23.1	0.6763006	PX	23	fusional/fissional/fusional	#N/A	#N/A
[283]	a.28.1	0.3526012	PX, SPX, PAX, SPAX	22, 22, 21, 21	fusional	ACP-like	cell periphery
[1556]	d.142.1 a.92.1 c.30.	0.6763006	PX	22	fusional/fissional/fusional	#N/A	#N/A

	1 d.142.1 c.24.1				signal		
[1085]	e.23.1 a.2 8.1 c.43.1 & e.23.1 a .28.1	0.6763006	PX	21	fusional/fi ssional/fu sional	#N/A	#N/A
{672}	a.4.1	0.6647399	PX, CX	8, 7	fissional/f usional	Homeodomain-like	regulation of epithelial cell differentiation involved in kidney development
{324}	c.73.1	0.3641618	PX	6	fissional/f usional	Carbamate kinase- like	heterocycle metabolic process
{460}	c.73.1 d.5 8.18& c.2. 1 d.81.1	0.5202312	CX	5	fusional/fi ssional	#N/A	#N/A
{734}	b.113.1 a. 156.1 g.3 9.1 c.37.1	0.6705202	CX	5	fusional/fi ssional	#N/A	#N/A
{13}	c.2.1 a.10 0.1	0.0693642	PX	5	fusional/fi ssional/fu sional	#N/A	#N/A
{58}	d.14.1	0.1791908	PX, SPX	5, 5	fusional	Ribosomal protein S5 domain 2-like	nucleic acid phosphodiester bond hydrolysis
{270}	g.39.1	0.3468208	PX	5	fissional/f usional	Glucocorticoid receptor-like (DNA-binding domain)	fibroblast growth factor receptor signaling pathway involved in ureteric bud formation
{388}	c.23.1	0.4046243	PX	5	fissional/f usional	CheY-like	regulation of multicellular organismal development

Figures

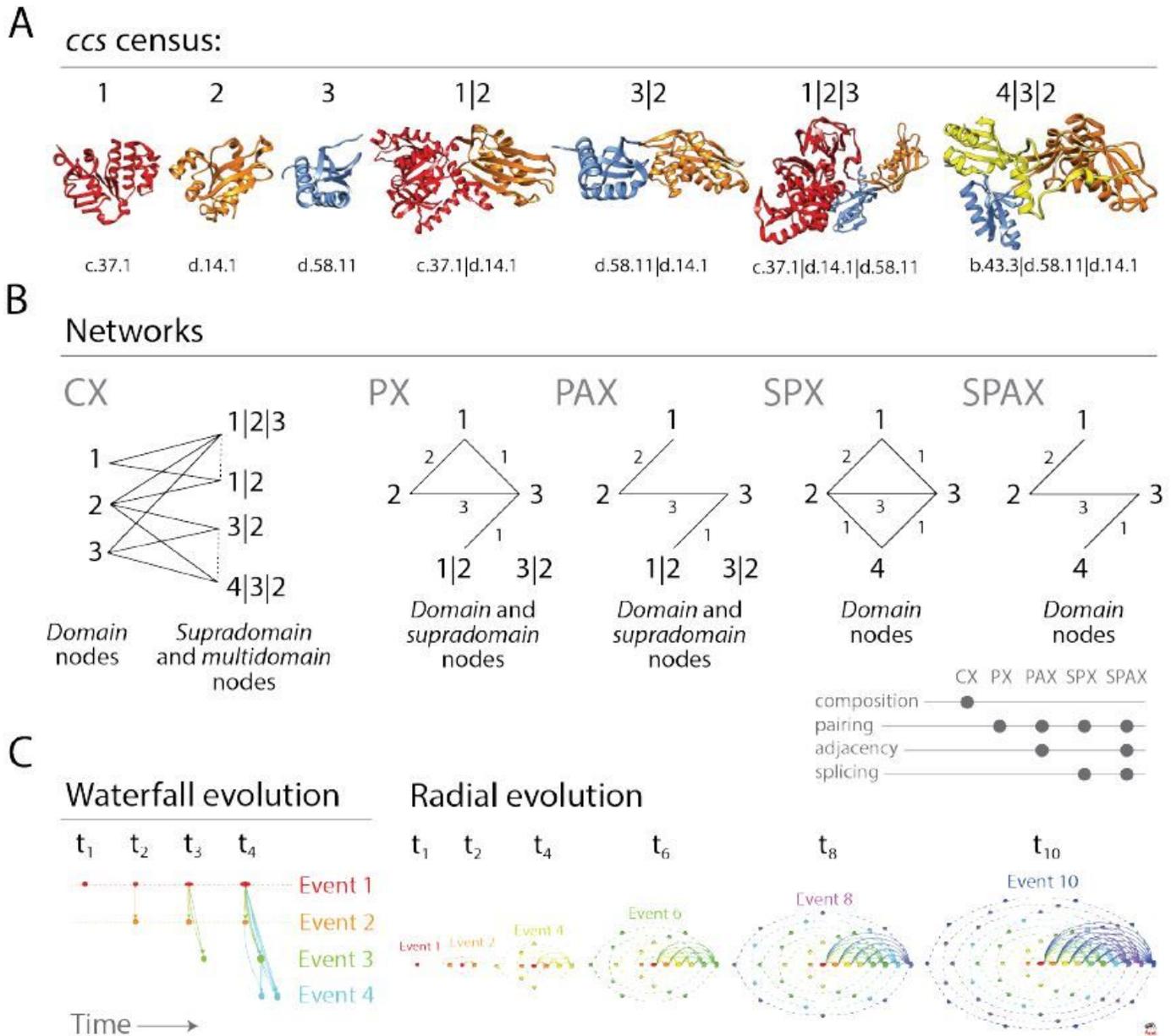


Figure 1

Networks of protein domain organization. (A) The genomic census of structural domains and their combinations defines SCOP concise classification string (*ccs*) descriptors of domains, supradomains and multidomains that are building blocks of networks. We illustrate the census with a sample from the entire entity set, comprising of 3 domains (1, 2 and 3), 2 supradomains (1|2 and 3|2) and 2 multidomains (1|2|3 and 4|3|2) that are common in dehydratase enzymes and elongation factors. *ccs* identifiers of structural domain constituents defined at fold superfamily (FSF) level are listed below the atomic models visualized in ribbon format with Chimera. (B) Five operative criteria for network generation capture the interactions among protein architecture nodes as networks grow in evolution. CX is a partial bipartite

network (projection-decomposable) that connects domain nodes to supradomain and multidomain nodes (which can connect to each other; hatched links) when present in multidomain proteins. PX connects domain and supradomain nodes when multidomain proteins are 'decomposed' into pairs of architectures, regardless of topological constraints. PAX borrows the PX criterion but respects topological constraints. SPX connects domain nodes spliced from architectures when domain pairs are present in proteins. SPAX connects domain nodes when adjacent domain pairs are present in proteins. (C) Chronological development of evolving networks. In 'waterfall evolution' layout, time progresses from left to right as 'discrete events' of network evolution progressively unfold the appearance of nodes and links (time-directed arrows known as arcs) from top to bottom, colored according to their age. Arc multiplicities describe link cardinality. Source-sink recruitments of architectures are visualized by horizontal and vertical elongations of node symbols, which describe their outdegree and indegree, respectively. As networks grow, the symbols of older nodes widen by outdegree accumulation, while those of younger nodes grow tall by indegree accumulation. In 'radial evolution' layout, the time-variant network grows by accumulating nodes in concentric rings (orbitals), each reflecting a time event. We illustrate radial evolution with 6 snapshots of a network growing to a size of 55 nodes as it unfolds from time t_1 to t_{10} . Nodes (n) in orbitals (r) grow at $r+1$ rate and only one node per orbital connects to single nodes in each of the other orbitals. Thus, outward links (o) of an orbital are $o=t-r-1$, where t is the current time. Inward links (i) of an orbital are $i=t-o-1=r$. Finally, total links of a network at any time are $t(t-1)/2$. The width and height of symbols represent the outdegree and indegree of nodes, respectively. Symbol sizes are shifted by 10 for a better visualization of nodes.

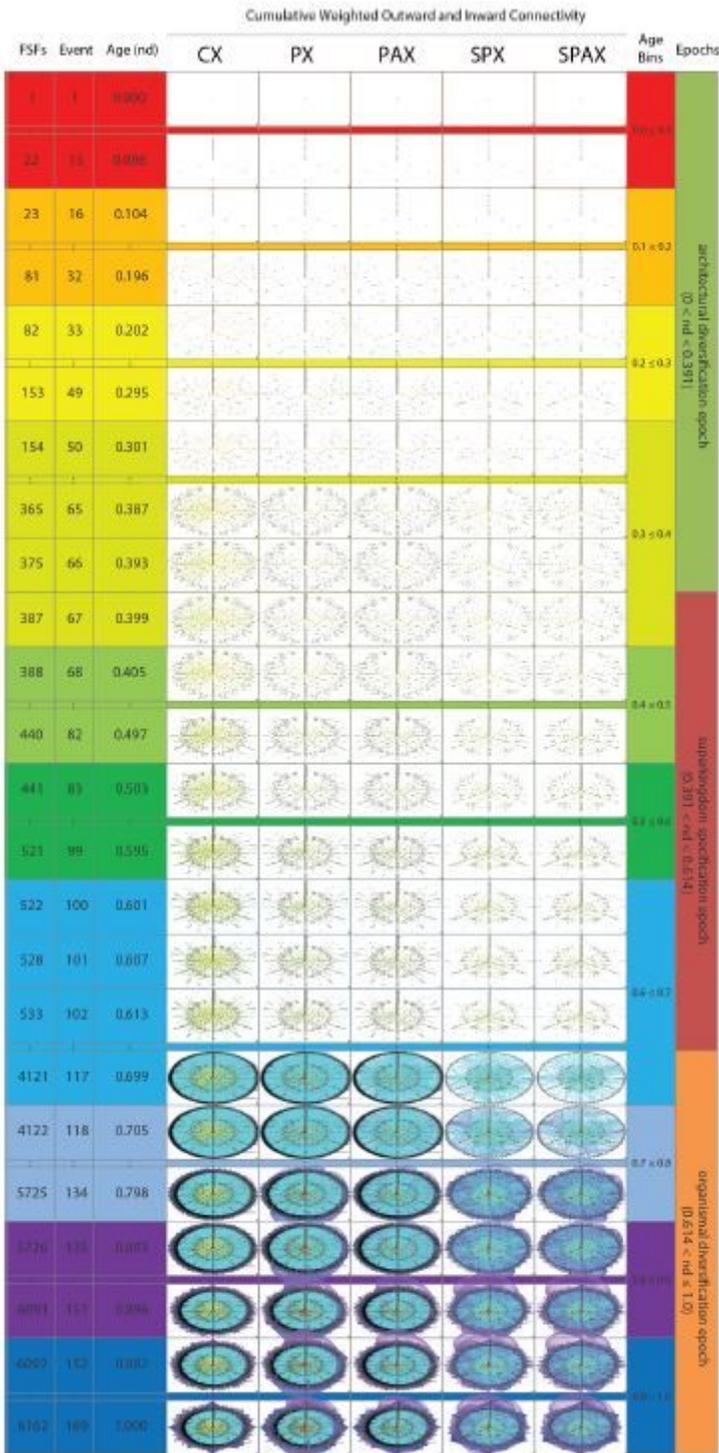


Figure 2

Evolving networks in radial evolution layout. Snapshots of network growth describe the evolution of 6,162 domain, supradomain and multidomain architectures or 1,643 domains spliced from them. They represent 24 out of 169 time events of the evolutionary timeline, which are indexed with evolutionary age (nd, ranging from 0.0 to 1.0), age bin (one of 10), and one of the 3 epochs of protein evolution (Wang et al., 2007). Age bins were custom RGB color coded to highlight the flow of time, from top to bottom. The evolving CX, PX, PAX, SPX and SPAX networks reveal the gradual evolutionary accumulation of nodes

and links. The sizes of the horizontal and vertical axes of the node symbols depict outward and inward weighted connectivity, respectively, with all weighted degree vectors shifted by 10 for visualization and inclusion of 0-degree nodes. The curved arcs describe recurring interactions between architectures that are accumulating along the successive events of the timeline. Arcs symbolize the flow of time from ancient to recent architectures and are color-coded according to the age of the more recent of the component nodes involved; arcs between contemporary nodes are excluded. Since, in pairwise networks the age of the most recent parent node could be assigned to the arc, the connectivity-defining pairing events are absent in the first (red) and the first and second (red, orange) bins of the PX and SPX and the PAX and SPAX networks, respectively. The angles of multiple arcs emerging from nodes are incremented by 2 to avoid overlap. Node RGB colors represent age. Grey-scale color of node borders depict fusional/fissional properties (Supplementary Fig. S3). Node shapes describe GO categories: circle, molecular function; squares, biological process; rhomboid, cellular component; triangle, unassigned.

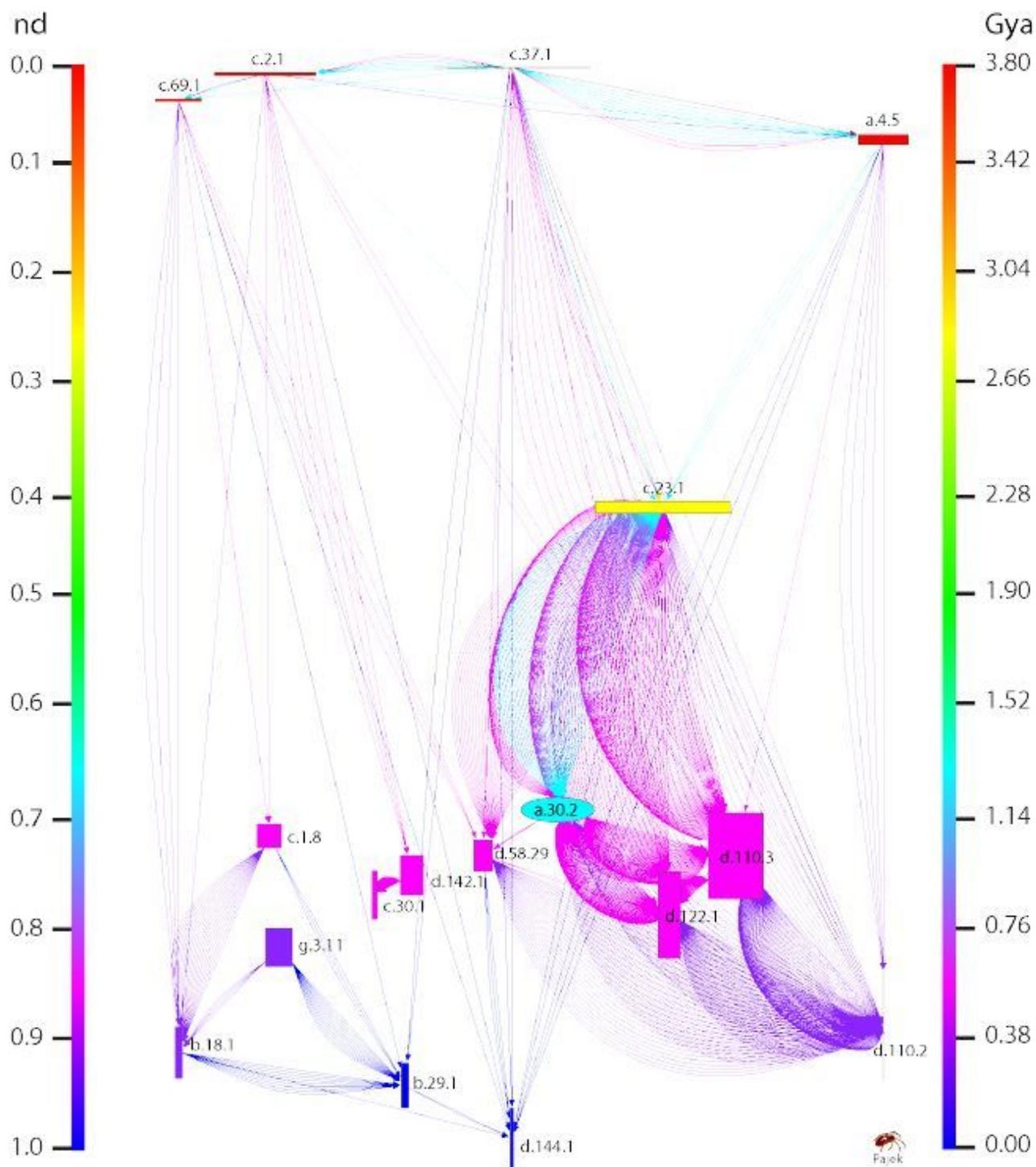


Figure 3

An extant SPX network in waterfall layout describing the evolution of spliced domains with the largest (100th percentile) network connectivity. The SPX network of 1,643 spliced domains was reduced with the restrictive criterion of excluding nodes with combined outdegrees and indegrees $\geq 99\%$ of those of the rest of the nodes. The set of arcs (arched arrows symbolizing flow of time) was also reduced to pairing events between domains in the 100th percentile connectivity and excluded those between contemporary

nodes. Nodes are arranged top-down and colored according to age (nd) on a relative 0-to-1 scale that describes evolutionary time events. Ages are also time-calibrated with a molecular clock of FSF domains, which uses fossils and microfossils, geochemical, biochemical, and biomarker data²⁰. FSF origin is given in billion years ago (Gya). Nodes were labeled with SCOP ccs domain descriptors. To showcase source-and-sink relationships, node symbol sizes were scaled proportional to the weighted outdegree and indegree along the horizontal and vertical axes, respectively. Weighted degrees were scaled as $\times 2 + 2$ to include 0-degree nodes for better visualization. The modular spread of nodes was based on VOS clustering (see methods). Arcs are color coded according to the age of the more recent of the component nodes involved; no arcs were present in the ancient most age bin (red) of the timeline. Angles of multiple arcs emerging from nodes are incremented by 2 to avoid overlap. See caption of Figure 2 for indexing of node colors and shapes.

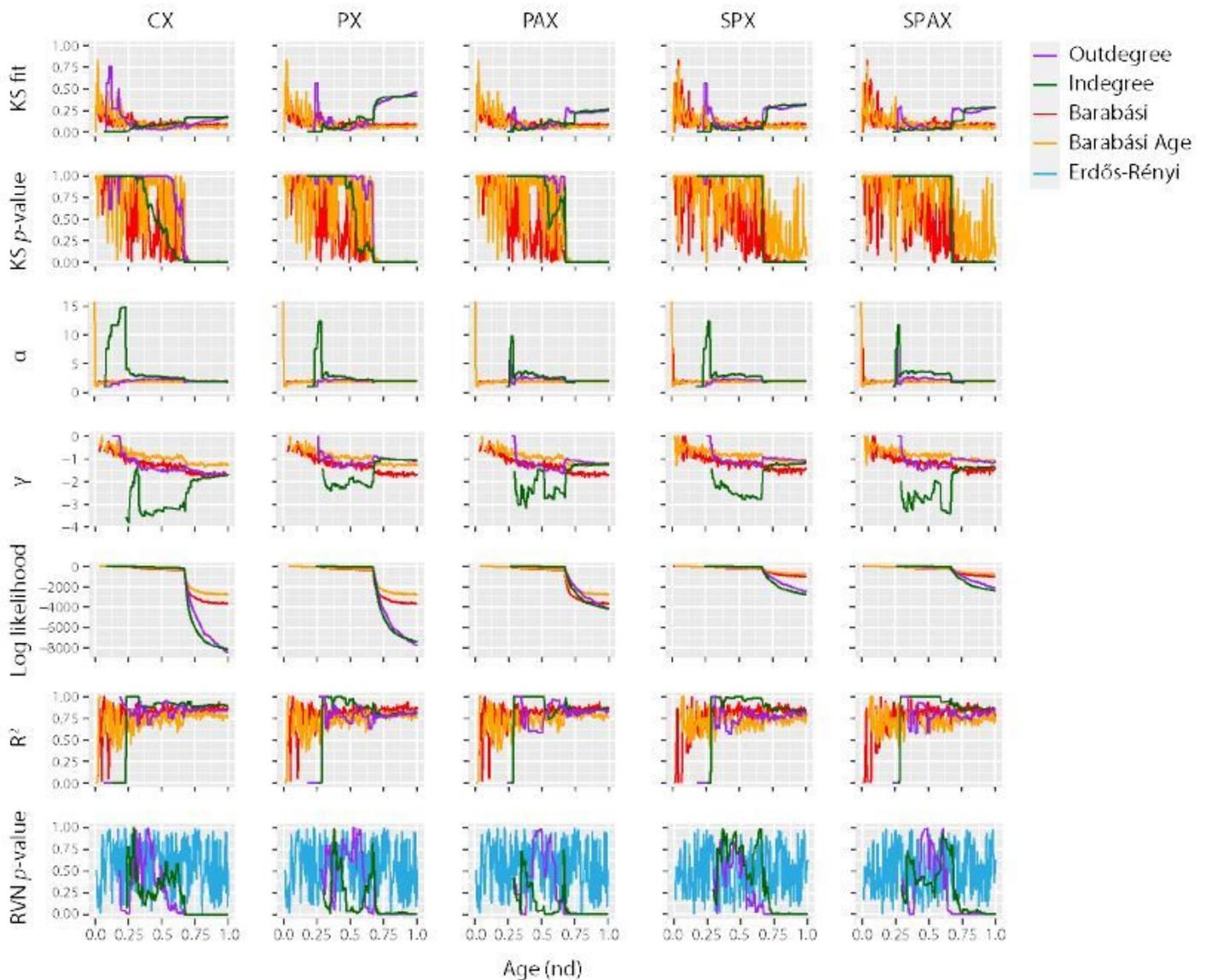


Figure 4

Statistical descriptors of power law and random behavior. Six indicators of preferential attachment were studied along the evolutionary timeline to explore processes of network growth, with network age (nd) indicated on a relative 0-to-1 scale. Outdegree and indegree connections were cumulative and weighted in evolving networks. Barabási (red) and Barabási-Age (orange) networks were included as control sets. The Barabási model specifies the probability of preference of an old node as $P_i \sim k_i^\alpha$ while the Barabási-Age model grants heavier power law properties to older nodes (exhibiting smaller nd) with $P_i \sim (k_i^\alpha)(l_i^\beta)$, where k_i is the indegree of node i of the current event, α is the preferential attachment exponent ($\alpha = 1$ for linear preferential attachment), l_i is the age of node i , i.e. the number of events elapsed since the node was added, with maximum number measured by the 'aging.bin' parameter, and β is the aging exponent ($\beta = 1$ for linear increases in probability of preference of an older node with high l_i). Power law indices include: (i) the KS fit statistic that compares the input degree data distribution with the fitted power law distribution (smaller scores denote better fit); (ii) the KS p-value, which rejects the null hypothesis that degree data was drawn from the fitted power-law distribution when less than $\alpha=0.05$; (iii) the exponent of the fitted power-law distribution (α); (iv) the slope of power-law linear regression model (γ); (v) the log-likelihood of the fitted parameters; and (vi) the coefficient of determination (R^2) that measures the percentage of degree data that fits the linear model. The randomness of the evolving networks was quantified by the p-value of an approximated beta distribution from the rank version of von Neumann's Ratio Test for Randomness⁶² (RVNp-value). The alternate hypothesis was non-randomness. Comparative graphs of strictly random Erdős–Rényi control networks of corresponding sizes at the given time-events were also plotted. Lower KS fit, higher KS p-value, higher α , lower γ and near-zero likelihood, given lower RVNp-value, support power law behavior.

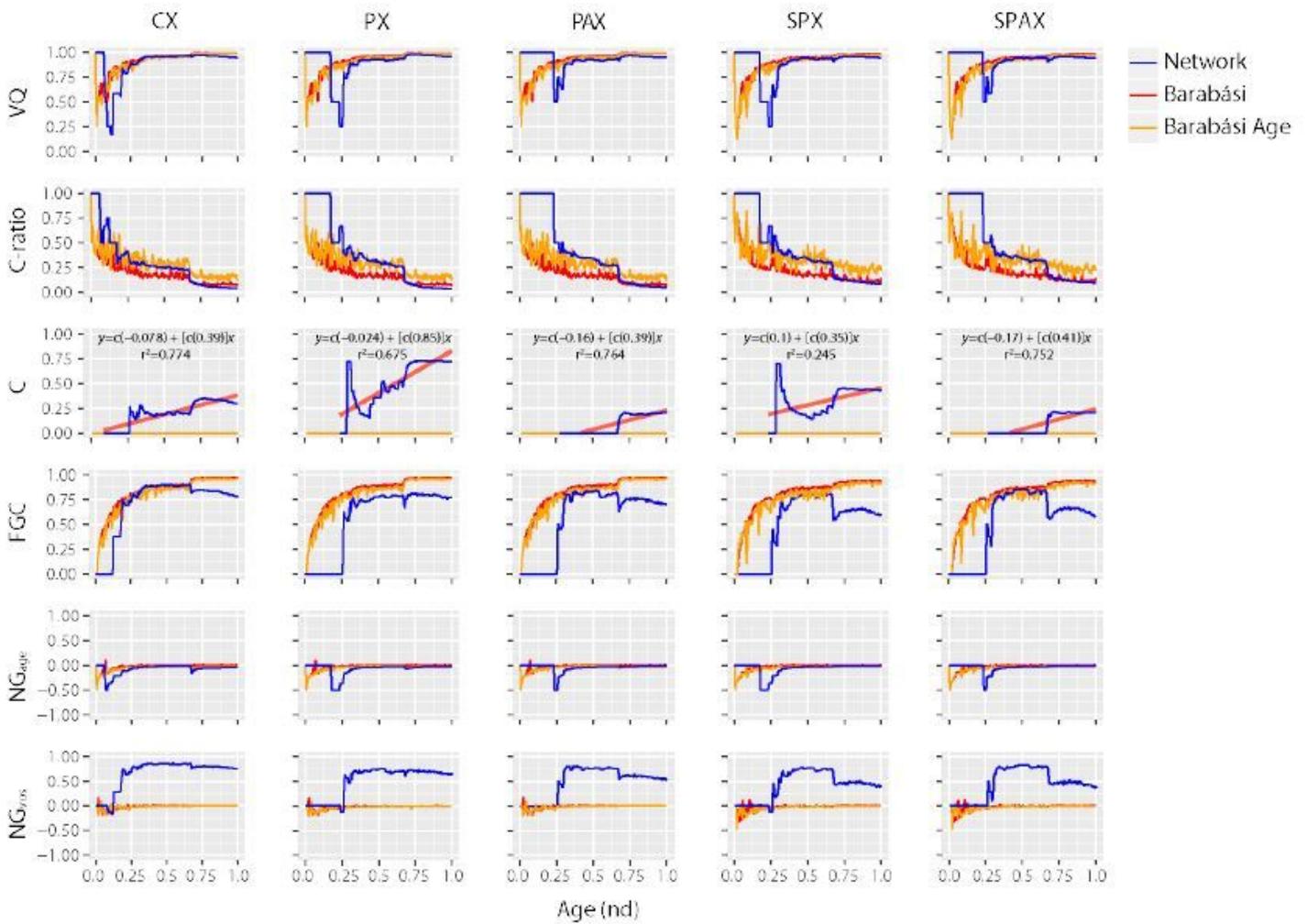


Figure 5

Network modularity. Six indicators of modularity were studied along the evolutionary timeline to explore the evolution of network structure, with network age (nd) indicated on a relative 0-to-1 scale. Modularity indices include the VOS Quality (VQ) index, the Clustering ratio (C-ratio), the average Clustering Coefficient (C), the Fast Greedy Community (FGC) index, and the Newman-Girvan index defined by age (NG_{age}) or VOS clustering (NG_{vos}). Modularity calculations required cumulative, undirected, and weighted connectivity input. The Barabási (red) and Barabási-Age (orange) models (see caption of Figure 4) were included as control sets. The regressions of C with age (nd) are shown as linear models (red lines) for each network together with supporting determination coefficients (R²).

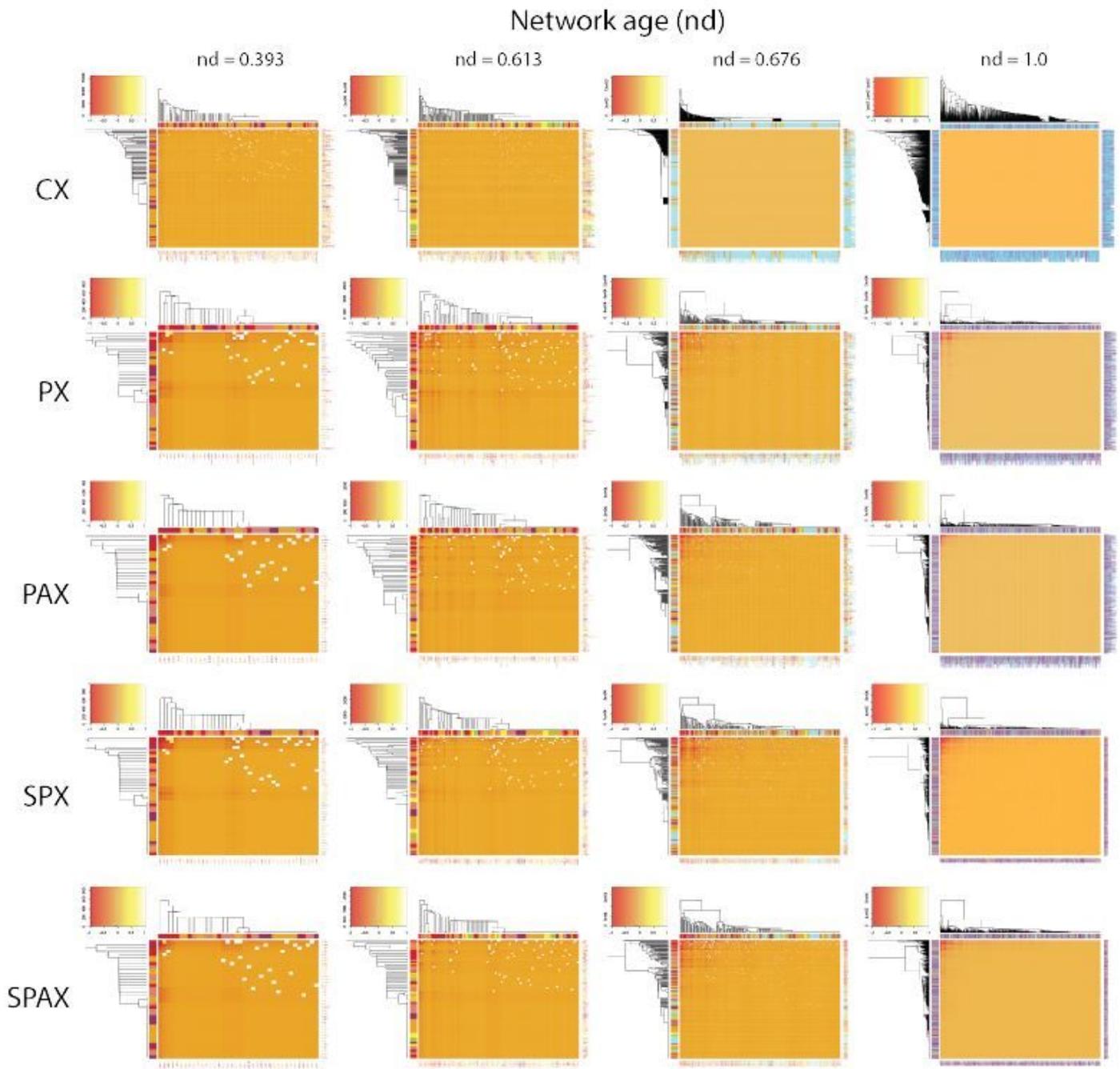


Figure 6

Evolution of modularity and hierarchical organization of networks over select events of the evolutionary timeline. NGage pairwise modularity values³⁹, scaled by \log_{10} of network-wide absolute modularity values, were used as input for the calculation of Euclidean distance matrices⁵⁴, which were visualized as heatmaps. Heatmap tiles represent modular strength between any two architectures relative to the respective strength of their linkages to other architectures of the network. The embedded dendrograms that define the order of rows and columns of the heatmaps were generated by hierarchical clustering of the distance matrices with the Ward's minimum variance method⁵⁵. The height of dendrograms represents dissimilarity between clusters while the clades show grouping rearrangements of

architectures. The top-left insets depict frequency histograms of the heatmap modularity values scaled from -1 to 1 (i.e. disassortative to assortative). The four panels describe growth of each evolving network (left-to-right). Network age corresponds to the middle approximate boundaries of the three evolutionary epochs of the protein world (Supplementary Fig. S2), i.e., end of 'architectural diversification' (nd = 0.393), end of 'superkingdom specification' (nd = 0.613), onset of the 'big bang' of domain organization at the start of 'organismal diversification' (nd = 0.676); and the present (nd = 1). Nodes were age-sorted ascendingly within clusters and labelled using standard SCOP nomenclature¹⁷. In the case of SPX and SPAX, nodes correspond to 1,643 domains mapped to the entity set of 6,162 architectures. The color-coding of bands and labels identifies the age of architectures (Supplementary Fig. S2). The relatively 'flatter' heatmap and 'skewed' dendrogram patterns of CX (typically at nd = 0.667 and nd = 1.000) are an artifact of unweighted distance matrices of CX, which contrast with the weighted ones of pairwise criterion-based networks. The most prominent clades correspond to the modules of the most ancient domain structures harboring the two major waves of architectural innovation. We also generated heatmaps of power-law control networks of corresponding sizes at the given time-events (Supplementary Fig. S10). When compared to the pairwise networks, the combined heatmap and dendrogram patterns of CX suggest a hidden switch from scale-freeness to modular behavior, eventually giving rise to hierarchical modularity with visible emergence of modules within modules.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [FigureS1.tif](#)
- [FigureS2.tif](#)
- [FigureS3.tif](#)
- [FigureS4.tif](#)
- [FigureS5.tif](#)
- [FigureS6.tif](#)
- [FigureS7.tif](#)
- [FigureS8.tif](#)
- [FigureS9.tif](#)
- [FigureS10.tif](#)
- [Video1SPXNetworkEvolution.mp4](#)
- [Video2SPXnd1PercentilesWaterfall.mp4](#)
- [Video3SPXHeatmapsCladsAnimation.mp4](#)
- [Manuscriptsupplement1.0.pdf](#)