

1 *A phylogenetic and proteomic reconstruction of eukaryotic chromatin evolution*

2

3 **Authors:** Xavier Grau-Bové^{1,2}, Cristina Navarrete^{1,2}, Cristina Chivas², Thomas Pribasni³,
4 Meritxell Antó⁴, Guifré Torruella⁵, Luis Javier Galindo⁵, Bernd Franz Lang⁶, David Moreira⁵,
5 Purificación López-García⁵, Iñaki Ruiz-Trillo^{4,7}, Christa Schleper³, Eduard Sabido^{1,2}, Arnau Sebé-
6 Pedrós^{1,2*}

7 1. Centre for Genomic Regulation (CRG), Barcelona Institute of Science and Technology (BIST), Barcelona
8 08003, Catalonia, Spain

9 2. Universitat Pompeu Fabra (UPF), Barcelona 08003, Catalonia, Spain

10 3. Department of Functional and Evolutionary Ecology, Archaea Biology Unit, University of Vienna, Djerassi-
11 platz 1, 1030 Vienna, Austria

12 4. Institut de Biologia Evolutiva (CSIC-Universitat Pompeu Fabra), Passeig Marítim de la Barceloneta, 37-49,
13 Barcelona 08003, Catalonia, Spain.

14 5. Unité d'Ecologie Systématique et Evolution, CNRS, Université Paris-Saclay, AgroParisTech, Orsay, France

15 6. Department of Biochemistry, Robert Cedergren Centre, Université de Montréal, Montréal, Quebec H3C 3J7,
16 Canada

17 7. ICREA, Pg. Lluís Companys 23, Barcelona 08010, Catalonia, Spain.

18

19 * Corresponding author: arnau.sebe@crg.eu

20

21

22 **Abstract**

23 Histones and associated chromatin proteins have essential functions in eukaryotic genome organiza-
24 tion and regulation. Despite this fundamental role in eukaryotic cell biology, we lack a
25 phylogenetically-comprehensive understanding of chromatin evolution. Here, we combine com-
26 parative proteomics and genomics analysis of chromatin in eukaryotes and archaea. Proteomics
27 uncovers the existence of histone post-translational modifications in Archaea. However, archaeal
28 histone modifications are scarce, in contrast with the highly conserved and abundant marks we
29 identify across eukaryotes. Phylogenetic analysis reveals that chromatin-associated catalytic func-
30 tions (e.g., methyltransferases) have pre-eukaryotic origins, whereas histone mark readers and
31 chaperones are eukaryotic innovations. We show that further chromatin evolution is characterized
32 by expansion of readers, including capture by transposable elements and viruses. Overall, our study
33 infers detailed evolutionary history of eukaryotic chromatin: from its archaeal roots, through the
34 emergence of nucleosome-based regulation in the eukaryotic ancestor, to the diversification of
35 chromatin regulators and their hijacking by genomic parasites.

36

37

38 Introduction

39 The access to genetic information in eukaryotes is controlled by a manifold nucleoproteic interface
40 called chromatin. This nucleosomal chromatin environment defines a repressive ground state for
41 transcription and other DNA-templated processes in eukaryotic genomes^{1,2}. Multiple components
42 associated with chromatin underlie elaborate eukaryotic genome regulation, allowing the differen-
43 tial access to genetic information in time/space and the maintenance of the resulting regulatory
44 states³⁻⁶. Moreover, chromatin-based regulation is essential in repressing parasitic genomic ele-
45 ments, like transposons and viruses⁷⁻¹¹.

46 The main protein components of eukaryotic chromatin are histones. All eukaryotes have four major
47 types of histones (H2A, H2B, H3 and H4), which are combined as an octamer to form the basic
48 repetitive unit of the chromatin: the nucleosome. Canonical histones are among the most highly
49 conserved proteins across eukaryotes¹² and, in addition, unique histone variants (paralogs of one of
50 the four major histone types) are found in many species, often associated with particular regulatory
51 states¹³⁻¹⁷. Histone chemical modifications, including acetylations and methylations play a central
52 role in genome regulation and transgenerational epigenetic inheritance^{3,18-21}. These chemical moi-
53 eties, known as histone post-translational modifications (hPTMs), are added and removed by specific
54 enzymes ('writers', e.g., histone methyltransferases or acetylases; and 'erasers', e.g., histone
55 demethylases and deacetylases). Some hPTMs (e.g., most acetylations) have a generic effect on
56 nucleosome stability, while others are bound by specific proteins or protein complexes. These are
57 often referred to as 'readers' and include proteins like HP1, which binds to H3K9me3, as well as a
58 myriad of other proteins encoding Chromo, PHD, Tudor and Bromo structural domains, among oth-
59 ers²²⁻²⁴. Finally, nucleosome remodellers (like SNF2 proteins) and histone chaperones are additional
60 important players in chromatin regulation, by mediating chromatin opening, nucleosomal assembly,
61 and histone variant interchanges²⁵⁻²⁸.

62 All eukaryotes studied to date possess histone-based chromatin organization, with the sole excep-
63 tion of dinoflagellates, which nonetheless encode for histone proteins in their genomes²⁹. Beyond
64 eukaryotes, histones have also been identified in Archaea, where they have been shown to form
65 nucleosomal structures³⁰⁻³³. However, unlike eukaryotic histones, the few archaeal histones experi-
66 mentally characterized so far (i) generally lack disordered N-terminal tails; (ii) do not have any
67 known post-translational modifications³⁴; and (iii) do not seem to impose a widespread, genome-
68 wide repressive transcriptional ground state^{33,35}. Thus, chromatin-based elaborate genome regula-
69 tion is often considered a eukaryotic innovation^{36,37}.

70 From a phylogenetic perspective, our understanding of chromatin components and processes de-
71 rives from a very small set of organisms, essentially animal, fungal and plant model species plus a
72 few parasitic unicellular eukaryotes. Additional efforts have sampled specific aspects of chromatin
73 regulation, such as histone modifications or their genome-wide distribution, in non-model animal
74 species^{38,39}, fungi (*Neurospora crassa* and *Fusarium graminearum*)^{40,41}, and five other eukaryotes:
75 the unicellular holozoan *Capsaspora owczarzaki*⁴², the dinoflagellate *Hematodinium* sp.²⁹, the
76 brown alga *Ectocarpus siliculosus*⁴³, the amoebozoan *Dictyostelium discoideum*⁴⁴, and the ciliate
77 *Tetrahymena thermophila*^{45,46}. However, these organisms represent a tiny fraction of eukaryotic

78 diversity. Hence, we lack a systematic understanding of the evolution of eukaryotic chromatin mod-
79 ifications and components⁴⁷.

80 In order to infer the origin and evolutionary diversification of eukaryotic chromatin, we performed a
81 joint comparative analysis of histone proteomics data from 30 different eukaryotic and archaeal
82 taxa, including new data for 23 species. In parallel, we analyzed the complement of chromatin-
83 associated gene families in an additional 172 eukaryotic genomes and transcriptomes. This compre-
84 hensive taxon sampling includes representatives of all major eukaryotic lineages, as well as multiple
85 free-living members of enigmatic early-branching eukaryotes (e.g., jakobids, malawimonads,
86 *Meteora* sp. and ancyromonads, as well as Collodictyonida, Rigifilida and Mantamonadida
87 (CRuMS); **Fig. 1a**). In addition, in order to trace the pre-eukaryotic origins of these chromatin gene
88 families, we systematically searched for orthologs in archaeal, bacterial and viral genomes. Specifi-
89 cally, we reconstructed the evolutionary history of enzymes involved in chromatin modification and
90 remodelling; as well as the conservation of the hPTMs effected by these enzymes. Our comparative
91 genomics and proteomics suggest a concurrent and early origin of canonical histones, a core of qua-
92 si-universal hPTMs, and their corresponding enzymatic effectors. We also identify independent ex-
93 pansions in hPTM reader gene families across eukaryotes and document evidence of the capture of
94 these reader domains by parasitic genomic elements. Overall, this work provides a
95 phylogenetically-informed framework to classify and compare chromatin components across the
96 eukaryotic tree of life, and to further investigate the evolution of hPTM-mediated genome regula-
97 tion.

98 **Results**

99 **Comparative proteomics of eukaryotic histone post-translational modifications.**

100 We analyzed the phylogenetic distribution and evolutionary history of histone proteins. To this end,
101 we surveyed the presence of histone-fold proteins across 172 eukaryotic and 4,226 archaeal taxa,
102 using HMM searches (**Fig. 1a,b and Supplementary Table 1**). Histone proteins are found in all
103 eukaryotic genomes. We clustered the identified 8,576 histone-encoding proteins using pairwise
104 local alignments and then classified individual sequences in these clusters based on pairwise align-
105 ments to a reference database⁴⁸ (**Fig. 1a and Supplementary Fig. 1a**). This reveals four broad clus-
106 ters corresponding to the four main eukaryotic histones (H2A, H2B, H3, and H4) and their variants
107 (H2A.Z, macroH2A, and cenH3), as well as a fifth cluster composed of archaeal HMfB homologs.
108 Finally, this classification also uncovers three large connected components composed of transcrip-
109 tion factors with histone-like DNA binding domains, which are widely distributed in eukaryotes
110 (POLE3, POLE4, DR1) and/or archaea (NFYB). Further analysis of the genomic distribution of
111 these histone genes shows a frequent occurrence of H3-H4 and H2A-H2B pairs in head-to-head
112 orientation (5' to 5'), strongly indicating co-regulation across eukaryotes (**Supplementary Fig. 1b,c**
113 **and Supplementary Table 2**).

114 Next, we investigated the distribution and conservation of hPTMs across major eukaryotic groups
115 and Archaea, including methylations, acetylations, crotonylations, phosphorylations, and
116 ubiquitylations. To this end, histones from 19 different eukaryotic species were extracted, chemical-
117 ly derivatized⁴⁹ and analyzed by mass-spectrometry (**Fig. 1c and Supplementary Table 3**), adding

118 to previously available hPTM proteomics data for additional seven species. Our extensive taxon
119 sampling covers all major eukaryotic groups, as well as hitherto unsampled early-diverging eukary-
120 otic lineages—such as the malawimonad *Gefionella okellyi*, the discoban *Naegleria gruberi*, or the
121 ancyromonad *Fabomonas tropica*—, thus providing a comprehensive comparative framework for
122 evolutionary inference.

123 We focused first on hPTMs present in canonical histones, as defined by their highly conserved *N*-
124 terminal regions, phylogenetic analyses, and sequence similarity to curated reference canonical his-
125 tones (**Fig. 1d**; see Methods). hPTMs are detected in all canonical histones from all species. After
126 correcting by sequence coverage, we observe that hPTMs are particularly abundant in H3 canonical
127 histones (median = 23.5 hPTMs per species, mean = 24.3), compared with H2A, H2B and H4 (me-
128 dians between 6.5 and 9, means between 9.5 and 13.4; **Supplementary Fig. 2a**). Holozoan canoni-
129 cal H2As (*Homo sapiens*, *Sycon ciliatum* and *Capsaspora owczarzaki*) represent an exception to
130 this trend and contain similar number of modifications to H3s in these species. We also examined
131 the reproducibility of hPTM detection across replicate samples, showing that the majority of hPTMs
132 (87.5%) can be found in more than one sample (**Supplementary Fig. 2b,c**). Despite this, it is worth
133 emphasizing that our data may contain false negatives, beyond the lack of coverage for particular
134 residues that we systematically report. For example, some marks might be globally too scarce in the
135 nucleosomes of a particular species, while other modifications like phosphorylations and
136 ubiquitination are difficult to detect by mass-spectrometry without dedicated peptide-enrichment
137 protocols.

138 Canonical H3 and H4 *N*-terminal tails contain the majority of phylogenetically-conserved hPTMs,
139 in stark contrast with the relative paucity of conserved hPTMs in canonical H2A and H2B. A strik-
140 ing example of paneukaryotic conservation comes from the acetylation of the H4 K5, K8, K12 and
141 K16 residues (**Fig. 1d**, second panel), all of which mark gene expression-permissive chromatin en-
142 vironments in multiple eukaryotic species²². A similar conservation pattern is observed in the acety-
143 lation of a group of *N*-terminal H3 lysines (K9, K14, K18, K23, K27) associated with similar func-
144 tions, while other H3 acetylations are only found in a few species (e.g., residues K4, K56 and K79).
145 While acetylations are highly conserved, only seven histone H3/H4 methylations are broadly con-
146 served across eukaryotic lineages: H3K4me1/2/3, H3K9me1/2/3, H3K27me1/2/3, H3K36me1/2/3,
147 H3K37me1/2/3 and, more sparsely, H3K79me1/2 and H4K20me1. Many of these broadly con-
148 served marks have conserved roles in demarcating active (e.g., H3K4me) and repressive chromatin
149 states (e.g., H3K9me and H3K27me)^{22,42,50}. The scarcity of conserved hPTMs in H2A and H2B
150 canonical histones can partially explained by their higher degree of sequence divergence (**Fig. 1e**),
151 which is reflected in many non-homologous lysine residues (**Fig. 1d**). But even among homologous
152 positions, we found little evidence of conservation, with the exception of H2A K5ac (associated to
153 active promoters⁵¹) and, in fewer species, methylation of H2A K5 and H2B K5. Finally, we were
154 also able to identify phosphorylations in serine and threonine residues and a few instances of
155 ubiquitylation. In general, these marks show more restricted phylogenetic distributions than lysine
156 acetylation or methylation, even in the tightly conserved H3 and H4 histones. We can identify con-
157 served phosphorylations in H2A T120 and S122, which are shared by most opisthokonts, and the
158 ubiquitylation of H2A K119 only in some holozoan species.

159 Mass-spectrometry analysis detected histone variants in all species included in our study, suggesting
160 that they are relatively abundant in the chromatin of these eukaryotes (**Fig. 1e**). Most of these vari-
161 ants are lineage-specific, with the exception of the paneukaryotic variants H2A.Z, H3/cenH3 and
162 H3.3; and the macroH2A variant found in holozoans and *Meteora* sp. (belonging to an orphan eu-
163 karyotic lineage). Interestingly, we find hPTMs in the vast majority of detected variants, both con-
164 served and lineage-specific, particularly acetylations and methylations (**Fig. 1e and Supplemen-**
165 **tary Fig. 2d**). Overall, our comparative proteomic analysis suggests the existence of a highly con-
166 served set of canonical hPTMs of ancestral eukaryotic origin in H3 and H4, which co-exists with
167 less conserved hPTMs in H2A, H2B, and lineage-specific modifications in variant histones.

168

169 **Archaeal histones and histone post-translational modifications**

170 In contrast with the paneukaryotic distribution of histones, sequence searches show that only a frac-
171 tion of archaeal genomes encode for histones (28.1% of the taxa here examined; **Fig. 2a**). Archaeal
172 histones exhibit a patchy phylogenetic distribution, similar to other gene families shared with eu-
173 karyotes⁵². Among others, histones are present in Euryarchaeota, the TACK superphylum and
174 Asgard archaea^{12,53-56}. Asgard are generally considered to be the closest known archaeal rela-
175 tives of eukaryotes^{57,58}, although this sister-group relationship has been challenged by some
176 studies⁵⁹. Our extended sampling revealed that Asgard archaea histones, particularly in the
177 Lokiarchaeota and Heimdallarchaeota clades⁵⁵, often have lysine-rich *N*-terminal tails in the manner
178 of eukaryotic histones (**Fig. 2a-c**). These Asgard histones appear to be conserved across multiple
179 taxa, albeit without direct sequence similarity compared to canonical eukaryotic histones (**Supple-**
180 **mentary Fig. 1d**). When compared against eukaryotic sequences classified in HistoneDB⁴⁸, these
181 archaeal histones clearly cluster in a separate group and are most similar to either eukaryotic H4 or,
182 to a lesser degree, H3 canonical histones, in line with previous findings^{12,55,60}.

183 To identify potential archaeal hPTMs, we performed proteomics analysis of histones in three
184 Euryarchaeota (the Methanobacteriota *Methanobrevibacter cuticularis* and the Halobacteriota
185 *Methanospirillum stamsii* and *Methanosarcina spelaei*) and one Thaumarchaeota species
186 (*Nitrososphaera viennensis*; **Fig. 2b**). Mass-spectrometry detects histone proteins in all of them: 2-4
187 in the euryarchaeotes (with 27-90% protein coverage) and one in the thaumarchaeote (80% protein
188 coverage), including homologs with *N*-terminal tails encoded by each of the three euryarchaeotes in
189 our survey (22-40 aa, 0.09-28 lysines per residue; **Fig. 2c**). Moreover, this proteomics analysis finds
190 evidence of hPTMs in archaeal histones. However, in comparison with eukaryotic histones, hPTMs
191 are extremely scarce in archaeal histones. Specifically, we identify no hPTMs in *N. viennensis* and
192 *M. spelaei* (one and two histones detected, respectively), three acetylations and one methylation in
193 *M. stamsii* (in three out of four histones detected), and one acetylation and two methylations in *M.*
194 *cuticularis* (in two out of four histones; **Fig. 2b**, top). Interestingly, we find conserved lysine resi-
195 dues with shared modifications in *M. stamsii* and *M. cuticularis* (methylation in K54 and acetyla-
196 tion in K57; **Fig. 2b**, bottom). This result indicates that highly-abundant hPTMs represent a eukary-
197 otic innovation, likely linked to dynamic nucleosomal regulation in eukaryotes but not in Archaea.

198

199 **Taxonomic distribution of chromatin-associated proteins**

200 hPTMs are deposited and removed by specific modifying enzymes (‘writers’ and ‘erasers’), while
201 ‘reader’ protein domains found in diverse proteins bind and recognize specific hPTMs. For exam-
202 ple, Bromo and Chromo domains bind acetylated and methylated lysine residues, respectively. In
203 addition, the control of histone loading/eviction from specific genomic *loci* is mediated by chroma-
204 tin remodellers, like SNF2 proteins²⁷, and histone chaperones²⁶. To date, the classification and evo-
205 lutionary analysis of this chromatin machinery has been based on biased, partial taxonomic sam-
206 plings and has not employed phylogenetic methods⁶¹ (with rare exceptions^{12,27}), often resulting in
207 inaccurate orthologous relationships and confounded classification and naming schemes.

208 We sought to obtain a systematic, phylogenetics-based classification of histone remodellers, chap-
209 erones, readers, and modifiers in order to understand the evolutionary history of eukaryotic chroma-
210 tin (**Fig. 3a**). To this end, we (i) compiled a taxa-rich dataset of 172 eukaryotic genomes and
211 transcriptomes, covering all major eukaryotic supergroups and devoting particular attention to ear-
212 ly-branching, non-parasitic lineages (**Supplementary Table 1**), as well as genomic data from 4,226
213 Archaea, 24,886 Bacteria and 185,579 viral taxa; (ii) defined a protein structural domain as a proxy
214 for each gene family (**Supplementary Table 4**) and retrieved all genes in these genomes that con-
215 tained these domains; and (iii) inferred accurate orthology groups from phylogenetic analyses of
216 each gene class (next section).

217 We examined the taxonomic distribution and abundance of the major gene classes (**Fig. 3b,c**). Many
218 domains with chromatin-associated functions in eukaryotes are also present in Archaea and Bacte-
219 ria, albeit with scattered phylogenetic distributions (**Fig. 3b and Supplementary Fig. 3a,b**). Fami-
220 lies with prokaryotic homologs include mostly catalytic gene classes (writer, eraser and remodeller
221 enzymes), whereas readers and histone chaperones are virtually absent from prokaryotes (**Fig. 3b**).
222 Histone fold-encoding genes constitute a case in point for this patchy distribution of chromatin pro-
223 teins in prokaryotes: they are present in most archaeal phyla, but are absent in about half of the
224 sampled genomes within each (**Fig. 3b**). Yet, there is a qualitative difference between the phyloge-
225 netic distribution of archaeal and bacterial chromatin-associated gene classes: whereas archaeal
226 histones tend to co-occur with chromatin-associated gene classes, the bacterial complement of writ-
227 ers and erasers is much less conserved and is uncorrelated with the extremely rare presence of his-
228 tone-like genes (**Fig. 3d**).

229 Within eukaryotes, most gene structural classes associated with chromatin functions are ubiquitous-
230 ly distributed across all lineages here surveyed, supporting an early eukaryotic origin for the core
231 chromatin machinery (**Fig. 3b and Supplementary Fig. 3d**). In fact, the total number of chromatin
232 writer, eraser and remodeller enzymes remains remarkably stable across eukaryotes (**Fig. 3e**). The
233 only exception is the marked increase in genes encoding reader domains observed in lineages exhib-
234 iting complex multicellularity: animals, streptophyte plants, and, to a lesser degree, phaeophyte
235 brown algae (Stramenopila). This occurs partially due to the addition of new gene classes (e.g.,
236 SAWADEE in the Plantae *s.l.* + Cryptista lineage, or ADD_DNMT3 in bilaterians and cnidarians),
237 but also via the expansion of ancient, widely-distributed reader gene classes (e.g., Tudor, PHD,
238 Chromo or Bromo domains). These taxonomic patterns indicate that chromatin modifying and re-
239 modelling catalytic activities originated in prokaryotes, while reader and chaperone structural do-
240 mains are eukaryotic innovations.

241

242 **Phylogenetics of chromatin modifiers and remodellers**

243 To gain detailed insights into the origin and evolution of chromatin gene families, we used phyloge-
244 netic analysis to define orthology groups from paneukaryotic gene trees. We surveyed 172 eukary-
245 otic species and defined a total of 1,713 gene families (orthogroups) encompassing 51,426 genes,
246 95% of which were conserved in two or more high-ranking taxonomic groups (as listed in **Fig. 1a**),
247 and which included 51,426 genes in total (**Supplementary Table 5**). We annotated each gene fami-
248 ly according to known members from eukaryotic model species. For simplicity, we use a human-
249 based naming scheme throughout the present manuscript (unless otherwise stated), but we also pro-
250 vide a dictionary of orthologs in three additional model species (*Arabidopsis thaliana*, *Saccharomy-
251 ces cerevisiae* and *Drosophila melanogaster*; see **Supplementary Table 5**). This phylogenetic clas-
252 sification scheme of eukaryotic chromatin gene families, as well as the sequences and associated
253 phylogenetic trees, can be explored and retrieved in an interactive database: [https://sebe-
254 lab.shinyapps.io/chromatin_evolution](https://sebe-lab.shinyapps.io/chromatin_evolution)

255 We first investigated the potential pre-eukaryotic origins of these gene families/orthogroups by
256 comparing their phylogenetic distance to prokaryotic sequences and to other eukaryotic orthogroups
257 (**Fig. 4a**). Most eukaryotic gene families are more closely related to other eukaryotes than to pro-
258 karyotic sequences, supporting the idea that writers, erasers, remodellers and readers diversified
259 within the eukaryotic lineage, as previously noted for histones¹². This analysis also reveals a sub-
260 stantial fraction of eukaryotic gene families with close orthogroups in Archaea and Bacteria, which
261 pinpoints components that were (i) inherited from a prokaryotic ancestor during eukaryogenesis; (ii)
262 laterally transferred between eukaryotes and prokaryotes at later stages; or (iii) a combination of
263 both phenomena. For example, we identified a well-supported sister-group relationship between the
264 eukaryotic SIRT7 deacetylase and a clade of Asgard archaea Sirtuin enzymes (Heimdallarchaeota
265 and Lokiarchaeota), a topology compatible with an archaeal origin or ancient transfers to/from
266 Asgard and eukaryotes⁶²; whereas SIRT6 appears nested within other eukaryotic sequences (**Fig.**
267 **4b**, left). Likewise, the KAT14 acetylase is more closely related to bacterial enzymes than to other
268 eukaryotic acetylases (**Fig. 4b**, right).

269 Next, we mapped the phylogenetic distribution of orthogroups in order to infer the origin and diver-
270 sification of individual chromatin gene families (**Fig. 4c and Supplementary Fig. 4a**). Using prob-
271 abilistic inference of ancestral gene content, we reconstruct a rich Last Eukaryotic Common Ances-
272 tor (LECA) complement of chromatin-associated gene families: 65 acetylases (amongst which 61
273 were conserved in at least two of the most deeply sampled eukaryotic early-branching lineages,
274 namely Amorphea, Diaphoretickes, and Discoba); 20 deacetylases (19 in these early-branching eu-
275 karyotic lineages); 59 methyltransferases (55); 42 demethylases (38); 33 remodellers (33); and 25
276 chaperones (18) (**Fig. 4c and Supplementary Table 5**). The subsequent evolution of these families
277 is characterized by relative stasis, with few new orthologous families emerging in later-branching
278 eukaryotic lineages. Notable exceptions include the origin of KAT5 deacetylases and KMT5B/C
279 SET methyltransferases in Opisthokonta; KAT8 and SIRT7 in Holozoa; and Viridiplantae-specific
280 deacetylases (homologs of *A. thaliana* HDA7 and HDA14 deacetylases) and SETs (*A. thaliana*
281 PTAC14); among others.

282 In spite of their broad distributions across eukaryotes, many chromatin modifier families exhibit
283 variation in their protein domain architectures, likely conferring them functional properties such as
284 distinct binding preferences (**Supplementary Fig. 4b**). For example, most CREBBP/EP300
285 acetylases consist of a catalytic HAT_KAT11 domain and two TAZ and ZZ zinc finger domains, but
286 different lineages have acquired different reader domains: an acetylation-reading Bromo domain in
287 holozoans and stramenopiles, PHD in plants and some stramenopiles, and no known reader domains
288 in other lineages (e.g., in the fungal orthologs of the *S. cerevisiae* protein RTT109). A similar pat-
289 tern is apparent in SET methyltransferase families sharing a core catalytic domain (SET) harboring
290 variable DNA- and chromatin-interacting domains – animal SETDB1/2 homologs have MBD do-
291 mains that bind CpG methylated DNA, while plants have SAD_SAR domains with the same func-
292 tion; and holozoan ASH1L homologs encode Bromo and BAH readers, whereas phaeophytes en-
293 code PHD domains (**Supplementary Fig. 4b**). Other architectures, however, are much more con-
294 served, as exemplified by the presence of Tudor-knot and MYST zinc finger domains in most KAT5
295 deacetylases; or the ubiquitous co-occurrence of Helicase-C and SNF2_N domains in most
296 remodellers (**Supplementary Fig. 4b**).

297 Specific examples of evolutionarily conserved chromatin gene families include the catalytic core
298 and the subunits of well-studied chromatin complexes⁶³ like PRC1 (RING1/AB, PCGF), PRC2
299 (EZH1/2, SUZ12, EED, RBBP4/7) and Trithorax/MLL (MLL1/2/3/4, WRD5, ASH2L, RBBP5,
300 DPY-30; **Fig. 4d,e**). However, when we compared the distribution of these complexes with the
301 hPTMs they are related to, we found a generally poor co-occurrence (**Fig. 4f-h**). For example, or-
302 ganisms like *Dictyostelium discoideum* and *Creolimax fragrantissima* lack EZH1/2 orthologs, but
303 we detected H3K27me3 in these species; while *Thecamonas trahens* and *Naegleria gruberi* lack
304 Dot1 orthologs but have H3K79me marks. A poor correlation is also observed between the occur-
305 rence of H3K9me and that of SUV39H1 orthologs. An exception to this pattern is the ubiquitous
306 distribution of H4K16ac and the acetylase family KAT5/8⁶⁴ (**Fig. 4h**). These patterns suggest that
307 the specificity between hPTMs and their writers might not be completely conserved across eukary-
308 otes, with distinct members of the same gene classes (e.g., methyltransferases) performing similar
309 roles. In this context, reading domains present in writing/erasing enzymes (directly in the same pro-
310 tein or as part of multi-protein complexes) are likely to play a major role in the re-purposing of
311 chromatin catalytic activities.

312

313 **Evolutionary expansion of chromatin readers**

314 Multiple protein structural domains have been involved in the recognition of hPTMs, such as
315 Bromo and PHD domains binding to acetylated lysines or Chromo, MBT and Tudor domains bind-
316 ing to methylated lysines^{23,24}. These are generally small domains and can be found both as stand-
317 alone proteins as well as in combination with other domains, often catalytic activities such as hPTM
318 writers, erasers and remodellers. Thus, they are central in the establishment of functional connec-
319 tions between chromatin states. To understand the contribution of these reading domains to the evo-
320 lutionary diversification of chromatin networks, we studied in detail the phylogeny and protein ar-
321 chitecture of reader domains across eukaryotes.

322 We quantified the co-occurrence frequency of reader and catalytic domains, finding (i) that most
323 reader domains are present in genes without writer, eraser or remodeller domains (87%, **Fig. 5a**);
324 and (ii) that most cases of reader-catalytic co-occurrence involve PHD, Chromo and Bromo do-
325 mains (**Supplementary Fig. 5a**). For example, the conserved architecture of the paneukaryotic
326 CHD3/4/5 re-modellers includes Chromo readers in most species and PHD domains specifically in
327 animals and plants (**Supplementary Fig. 4b**). Likewise, PHD domains are often present in the
328 KMT2A/B and KMT2C/D SET methyltransferase; and the ASH1L family has recruited Bromo and
329 BAH domains in holozoans, and PHD in multicellular stramenopiles (**Supplementary Fig. 4b**). In
330 spite of these redundancies, reader families typically have independent evolutionary histories, as
331 illustrated by the fact that most reader domain-containing genes encode only one such domain
332 (92%, **Supplementary Fig. 5b**).

333 We next performed phylogenetic analyses of individual reader domains and reconstructed the gains
334 and losses of these reader gene families/orthogroups (**Fig. 5a**). Compared to the relative stasis of
335 catalytic enzyme families, this reader-centric analysis revealed a strikingly different evolutionary
336 pattern of lineage-specific bursts of innovation, particularly amongst PHD, Chromo and Bromo
337 genes, as well as Tudor in animals (**Fig. 5a and Supplementary Fig. 5c**). PHD, Chromo and
338 Bromo families also appeared as the most abundant in the reconstructed LECA reader domain rep-
339 ertoire, which amounted to 89 gene families (**Fig. 5a**, left). The distribution of gene family ages in
340 extant species also corroborates that more readers have emerged at evolutionarily more recent nodes
341 of the tree of life than catalytic gene families (**Fig. 5b**).

342

343 **Co-option of the chromatin machinery by transposable elements**

344 Further examination of the domain co-occurrence networks of readers revealed that Chromo and
345 PHD domains are often present together with protein domains found in transposable elements (TEs;
346 **Fig. 5c and Supplementary Table 6**), including retrotransposons (e.g., retrotranscriptases and
347 integrases; orange modules in **Fig. 5c**) and DNA transposons (e.g., DNA binding domains and
348 transposases; red modules). It is known that some TEs show insertion-preferences associated to
349 specific chromatin states⁶⁵, often mediated by direct chromatin tethering mechanisms⁶⁶. For exam-
350 ple, the Chromo domain of the MAGGY gypsy retrotransposon of the fungus *Magnaporthe grisea*
351 targets H3K9me regions⁶⁷. Reciprocally, some protein domains of TE origin, often DNA-binding
352 domains, have been co-opted into chromatin and transcriptional regulators⁶⁸. Thus, we decided to
353 explore in detail the occurrence of chromatin-associated domain (readers, but also catalytic do-
354 mains) linked to TEs in the 172 eukaryotic genomes in our dataset (**Fig. 5d**). Moreover, we used
355 available RNA-seq datasets in many of these species to validate some of these TE fusions (**Fig. 5d-**
356 **e**). A fully validated fusion gene would (i) come from a non-discontinuous gene model in the origi-
357 nal assembly, and (ii) have evidence of expression, with reads mapping along the entire region be-
358 tween the TE-associated domain and the chromatin-associated domain (**Supplementary Fig. S6**).

359 We identified 823 predicted gene models containing both chromatin- and TE-associated domains
360 (**Fig. 5d**). Whilst these TE fusions were not exclusive of reader domains, most such fusions in-
361 volved PHD and Chromo-encoding genes; followed by SNF2_N remodellers, SET
362 methyltransferases, and others. An homology search against a database of eukaryotic TEs revealed

363 that most of these candidate TE fusions could be aligned to known retrotransposons or DNA trans-
364 posons. For example, by way of validation, our analysis identifies the SETMAR human gene, a
365 previously-described fusion between a SET methyltransferase and a Mariner-class DNA transpos-
366 on⁶⁹. Overall, 31% of the candidate fusion genes were supported by valid gene models according to
367 our stringent criteria (**Fig. 5d**). Interestingly, we find very few cases of hypothetical fusions be-
368 tween TEs and Bromo domains, which recognize K acetylations and are otherwise highly abundant
369 across eukaryotes, and none of them is validated by RNA-seq data. This could be explained by the
370 detrimental effect of targeting TE insertions to sites of active chromatin demarcated by histone
371 acetylations, such as promoter and enhancer elements.

372 Some of these validated fusions have a broad phylogenetic distribution (**Fig. 5e**), such as a Gypsy-
373 ERV retrotransposon with a C-terminal Chromo domain (Unk. Chromo 2.1 in **Fig. 5e**) that is widely
374 distributed in animals and various microbial eukaryotes, and contains dozens of paralogs in verte-
375 brate *Danio rerio* or the charophyte *Chara braunii*, many of which are expressed. Another wide-
376 spread Gypsy-ERV retrotransposon with a Chromo domain is present in multiple expressed and
377 highly similar copies in the fungus *Rhizopus delemar* (**Fig. 5f,e**), suggesting a successful coloniza-
378 tion of this genome by this TE. By contrast, other TE fusions are taxonomically restricted to one or
379 few related species, such as the fusion of hAT activator DNA transposons with Chromo CBX and
380 CDY readers in the sponge *Ephydatia muelleri*; or multiple instances of fusions with Chromo and
381 PHD readers in cnidarians. A common fusion in cnidarians involves different retrotransposon clas-
382 ses with PHD domains orthologous to the PYGO1/2 protein (**Fig. 5e**), which is known to recognize
383 specifically H3K4me⁷⁰. Globally, this analysis reveals that recruitment of chromatin reading and
384 even modifying domains by TE has occurred in many eukaryotic species, in a way that might facili-
385 tate the evasion from suppressing mechanisms in the host genomes as suggested by the expansion
386 of Chromo-fused TEs in the genomes of *Chara braunii* (Viridiplantae), *Chromera velia* (Alveolata)
387 and *Rhizopus delemar* (Fungi).

388

389 **Chromatin components in viral genomes**

390 In addition to TEs, chromatin is also involved in the suppression of another type of genomic para-
391 sites: viruses. Some chromatin-related genes, including histones, have been found in viral genomes,
392 especially among the nucleocytoplasmic large DNA viruses – also known as giant viruses. Eukary-
393 otic core histones have been even hypothesized to have evolved from giant virus homologs, after
394 the discovery that certain Marseilleviridae genomes encoded deeply-diverging orthologs of the four
395 canonical histones⁷¹. These viral histones have been recently shown to form nucleosome-like parti-
396 cles that package viral DNA^{72,73}.

397 We analyzed the distribution and abundance of chromatin-related protein domains among viruses,
398 including data from 1,816 giant virus genomes. Based on structural domain searches, we identified
399 2,163 viral chromatin-related proteins (**Fig. 5g and Supplementary Table 6**). The majority of these
400 proteins are encoded by giant viruses (55%), followed by Caudovirales (37%). Among these two
401 groups, only giant virus genomes encode histones – specifically, the Iridoviridae, Marseilleviridae,
402 Mimiviridae, Pithoviridae, and Phycodnaviridae families. Concordantly with previous studies⁷⁴, we

403 also identify remodellers in all giant virus families; as well as less abundant components of the
404 chromatin writer/eraser/reader toolkit (**Fig. 5g**).

405 We then investigated the phylogenetic affinities of these viral chromatin proteins, starting with his-
406 tones (**Fig. 5h**). Our analysis recovers the phylogenetic affinity of Marseilleviridae histones with
407 specific eukaryotic histone families⁷¹, and makes this pattern extensive to Mimiviridae, Iridoviridae,
408 and Pithoviridae giant viruses (**Fig. 5h**), with the caveat of the ambiguous clustering of the H4-like
409 viral histones with either H4 eukaryotic or archaeal HMfB genes. In all these lineages, we identify
410 genes encoding two histone-fold domains orthologous to H2B + H2A (inset table in **Fig. 5h**),
411 whereas the H4 + H3 histone doublet genes appears to be exclusive to Marseilleviridae. By contrast,
412 histone homologs in Phycodnaviridae, Pandoraviridae (also giant viruses), and Polydnviridae
413 (*incertae sedis*) are never found as either doublets or as early-branching homologs of eukaryotic
414 histones, suggesting recent acquisition from eukaryotes.

415 Unlike histones, most of the viral chromatin-associated genes exhibited a mixture of prokaryotic
416 and eukaryotic phylogenetic affinities and often lack affinity to any specific eukaryotic gene family
417 (**Fig. 5i and Supplementary Fig. 7**). Viral readers, on the other hand, are often embedded within
418 eukaryotic clades in gene trees and are similar to *bona fide* eukaryotic families, exhibiting topolo-
419 gies consistent with recent, secondary acquisitions. This is the case of BIRC2/3/XIAP readers wide-
420 spread in the Baculoviridae, which encode BIR domains that are often hijacked from their hosts⁷⁵.
421 We also find a number of viral Chromo-encoding genes, which fall in two main taxonomic catego-
422 ries: (i) giant virus homologs of the eukaryotic CBX1/3/5 family (present in Mimiviridae,
423 Iridoviridae and Phycodnaviridae); and (ii) homologs from various Adintoviridae, which are closely
424 related to animal Chromo genes encoding *rve* integrase domains⁷⁶ (**Fig. 5i**). Finally, we also identify
425 a handful of eukaryotic-like viral genes with deep-branching positions relative to core eukaryotic
426 gene families, as seen in histones (**Fig. 5h**). This includes Mimiviridae homologs of the eukaryotic
427 methyltransferases SMYD1-5 and DOT1 (**Supplementary Fig. 7d,e**), as well as SNF remodeller
428 families with homologs in distinct giant virus clades (HLTF/TTF2 in Phycodnaviridae, Mimiviridae
429 and Iridoviridae). These results indicate that cases of horizontal transfer from eukaryotes to viruses
430 are common in different chromatin-related gene families, including histones. Therefore, it is likely
431 that basally-branching giant virus histones were similarly acquired from a stem eukaryotic lineage
432 and this would explain the observed histone tree topology with extant eukaryotic species. In any
433 case, most of the eukaryotic chromatin machinery appears to have cellular roots.

434 Discussion

435 Our comparative proteogenomics study reconstructs in detail the origin and evolutionary diversifi-
436 cation of eukaryotic chromatin components, from post-translational modifications to gene family
437 domain architectures. We looked first at the pre-eukaryotic roots of chromatin. Multiple aspects of
438 archaeal chromatin have been studied in recent years, including nucleosomal patterns³¹ and the
439 structure of the archaeal nucleosome³⁰. A recent taxonomic survey of archaeal nucleoid-associated
440 proteins revealed multiple independent diversifications of DNA-wrapping proteins and a strong
441 association between high levels of chromatinization and growth temperature, overall suggesting a
442 structural, non-regulatory role for archaeal chromatin⁷⁷. Our proteomics data support this notion by
443 showing the scarcity of hPTMs in four species belonging to two different archaeal lineages

444 (Euryarchaeota and Thaumarchaeota). An earlier proteomics study reported the complete absence of
445 hPTMs in the euryarchaeote *Methanococcus jannaschii*³⁴. Here we do identify a few instances of
446 modified lysine residues in Euryarchaeota, which is in line with the recently reported acetylations in
447 *Thermococcus gammatolerans* histones⁷⁸. It remains to be seen if hPTMs are frequently present in
448 Asgard and other unsampled archaeal lineages, where other eukaryotic-like features have been
449 found^{57,79,80}. In fact, some of these Asgard, particularly Lokiarchaeota, encode for histones with
450 long, K-rich N-terminal tails but that bear no similarity with eukaryotic histones and are, therefore,
451 most probably the result of convergent evolution. Interestingly, Lokiarchaeota genomes also fre-
452 quently encode histone modifiers such as SET methyltransferases and MOZ_SAS acetylases. How-
453 ever, overall our results suggest that extensive usage of hPTMs is an eukaryotic innovation (**Fig.**
454 **6a**). Similarly, while we find the majority of catalytic domains of hPTM writers, hPTM erasers and
455 chromatin remodellers in Archaea and even Bacteria, these appear only scattered in a small fraction
456 of the examined taxa. In contrast, hPTM reader domains and histone chaperones are eukaryotic in-
457 novations, further supporting the idea that the functional readout of hPTMs and the role for histone
458 variants in defining chromatin states are both exclusive to eukaryotes (**Fig. 6a**).

459 The origin of eukaryotes represents a major evolutionary transition in the history of life⁸¹. Thanks to
460 sequencing and comparative analysis of archaeal and eukaryotic genomes, we also have a detailed
461 reconstruction of the massive innovation in gene repertoires that occurred at the origin of eukary-
462 otes. This gene innovation in the Last Eukaryotic Common Ancestor (LECA) includes cytoskeletal
463 proteins and associated motors like myosins^{82,83} and kinesins⁸⁴, vesicle trafficking apparatus⁸⁵,
464 splicing machinery⁸⁶, ubiquitin signalling systems⁸⁷ and a large repertoire of sequence-specific tran-
465 scription factors³⁷. Combining parsimony analysis and knowledge on gene function in extant line-
466 ages (mostly vertebrates, yeast and plants), our results allow us to reconstruct a complex LECA
467 repertoire of hPTMs and associated writing, eraser and reader gene families (**Fig. 6b,c**). We infer 23
468 to 29 highly-conserved lysine acetylations in canonical histones (e.g., H3K9ac and H3K27ac) and a
469 repertoire of 65 and 20 histone acetylase and deacetylase families, respectively. With the exception
470 of H4K16ac⁶⁴, most histone acetylations are thought to exert a generic, perhaps additive, effect on
471 the opening of chromatin²². As such, acetylation marks like H3K27ac have been found to be en-
472 riched in promoters of active genes in diverse eukaryotes⁴². In contrast, histone methylations often
473 have very specific readouts and they can be linked both to active and repressive chromatin states.
474 We infer between 13 and 25 conserved methylated lysine residues in LECA histones, including
475 marks typically associated to active promoters (H3K4me1/me2/me3), gene bodies (H3K36me3,
476 H3K79me1/2, H4K20me1), and repressive chromatin states (H3K9me2/me3, H3K27me3,
477 H4K20me3)^{88,89}. Finally, we also infer the existence of five histone variants in the LECA (cenH3,
478 H3.3, H2A.Z, macroH2A and H2A.X), as well 33 chromatin remodellers (e.g., EP400/SWR1 and
479 INO80, involved in loading and removal of H2A.Z, respectively) and 25 histone chaperones (e.g.,
480 ASF1A/B and NPM1/2/3). This indicates that, in addition to an extensive repertoire of hPTMs, the
481 regulation of nucleosomal histone composition was also an important feature in the LECA.

482 Chromatin evolution after the origin of eukaryotes is characterized by an expansion of lineage-
483 specific histone variants harboring unique hPTMs and a net expansion in the number of reader gene
484 families, as opposed to the relatively static catalytic gene families (writers, erasers and remodellers).
485 This is particularly relevant as it suggests extensive remodelling of chromatin networks during eu-

486 karyote evolution, that is, changes in the coupling of particular hPTMs to specific functional chromatin states. An example of such changing state-definitions comes from looking at the hPTMs associated to TEs in different organisms: H3K9me3+H4K20me3 in animals, H3K27me3 in some plants⁹⁰, H3K79me2+H4K20me3 in the brown multicellular algae *Ectocarpus siliculosus*⁴³, and H3K9me3+H3K27me3 in the ciliate *Paramecium tetraurelia*⁹¹. In the context of the histone code hypothesis^{3,20,92-94}, our findings indicate that, while there is an ancient core of conserved hPTMs across eukaryotes, evidence for a universal code/functional-readout is limited, with perhaps the exception of the highly conserved configuration of ancient hPTMs around active promoters across many eukaryotes⁴². Another interesting observation related to the evolution of chromatin networks is the capture of chromatin reader domains by TEs. We find evidence of this phenomenon in a number of species with a scattered phylogenetic distribution, suggesting that it is a recurrent process and that it often leads to the successful propagation of the TE in the host genome. We hypothesize that this process facilitates the targeting of TEs to specific chromatin states, as it has been described in the case of MBD DNA methylation readers captured by TEs^{95,96}.

500 In the future, a broader phylogenetic understanding of the genome-wide distribution of hPTMs, as well as the direct interrogation of hPTM binders in different species⁹⁷⁻⁹⁹, will be crucial to further clarify questions such as the ancestral role of specific hPTM and the co-option of ancient hPTMs into novel functions.

504

505 **Acknowledgements**

506 We want to thank Alex de Mendoza for critical input on the analysis of transposable element fusions. We also want to thank Josep Casacuberta for *Physcomitrella patens* samples, Harold J. G. Meijer for *Phytophthora infestans* samples, Maja Adamska for *Sycon ciliatum* samples, and Alistar Simpson for access to the *Gefionella okellyi* culture (made possible by his funding from NSERC, Canada). Research in A.S.-P. group was supported by the European Research Council (ERC) under the European Union's Horizon 2020 Research and Innovation Programme Grant Agreement (851647), the Spanish Ministry of Science and Innovation (PGC2018-098210-A-I00), the Centro de Excelencia Severo Ochoa scheme (SEV-2016-0571), and the Agencia Estatal de Investigacion. C.N. is supported by an FPI PhD fellowship from the Spanish Ministry of Economy, Industry and Competitiveness (MEIC). X.G.-B. is supported by a Juan de la Cierva fellowship (FJC2018-036282-I) from the Spanish Ministry of Economy, Industry and Competitiveness (MEIC). I.R.-T. was supported by a European Research Council Grant (616960). B.F.L. was supported by the Natural Sciences and Engineering Research Council of Canada (NSERC; RGPIN- 2017-05411) and by the 'Fonds de Recherche Nature et Technologie', Quebec. P.L.-G. and D.M. were supported by a Moore and Simons foundations grant (GBMF9739) and by European Research Council Advanced Grants (322669, 787904). Research in C.S. group was supported by the European Research Council (ERC) through project TACKLE (AdvGrant No. 695192).

522

523 **Author contributions**

524 A.S.-P. conceived the project. X.G.-B., C.C., I.R.T., C.S., E.S. and A.S.-P. designed experiments and analytical strategies. C. N., T.P., M. A. and A.S.-P. performed experiments. X.G.-B., C.C., and A.S.-P analyzed the data. T.P., G.T., L.J.G., D.M.,P.L.-G. and B.F.L. provided biological samples/cultures and genomic data. All

527 authors contributed to data interpretation. X.G.-B. and A.S.-P. wrote the manuscript with input from all au-
528 thors.

529

530 **Declaration of interests**

531 The authors declare no competing interests.

532

533

534 **MAIN FIGURES TITLES AND LEGENDS**

535 **Figure 1. Diversity of post-translational modifications in eukaryotic canonical and variant**
536 **histones. a**, Eukaryotic taxon sampling used in this study. Colored dots indicate the number of spe-
537 cies used in the comparative histone proteomics reconstruction, with solid dots indicating new spe-
538 cies added in this analysis. Numbers in brackets indicate the number of genomes/transcriptomes
539 used in the comparative genomics analyses. Dashed lines indicate uncertain phylogenetic relation-
540 ships. Complete list of sampled species in Supplementary Table 1. **b**, Networks of pairwise protein
541 similarity between histone protein domains in eukaryotes, archaea and viruses. Each node repre-
542 sents one histone domain, colored according to their best alignment in the HistoneDB database (see
543 Methods). Edges represent local alignments (bitscore ≥ 20). **c**, Schematic representation of the
544 hPTM proteomics strategy employed in this study. **d**, Conservation of hPTMs in eukaryotic his-
545 tones. hPTM coordinates are reported according to the amino-acid position in human orthologs (if
546 conserved). In H2A and H2B, question marks indicate the presence of hPTMs in stretches of lysine
547 residues of uncertain homology. In species with previously reported hPTMs, we further indicate
548 which variants were also identified in our reanalysis. Only positions with hPTMs conserved in more
549 than one species are reported (full table and consensus alignments available in Supplementary Table
550 3). **e**, Maximum likelihood phylogenetic trees of the connected components in panel b, correspond-
551 ing to eukaryotic histones (H3, H4, H2A, H2B). Canonical histones included in panel d and variant
552 histones detected are highlighted in red. hPTMs detected in non-canonical histones are indicated.
553 Bottom, distributions of pairwise phylogenetic distances between all proteins in each gene tree. Vio-
554 lin plots above each distribution represent the distribution of distances between reference histones
555 present in the HistoneDB database and histones with proteomic evidence included in our study, for
556 each of the main canonical (H3, H4, H2A, and H2B) and variant histones (H2A.Z and macroH2A).

557 **Figure 2. Archaeal histone diversity and post-translational modifications. a**, Distribution of
558 histones (fraction of taxa in each lineage) and histone tails (presence/absence) across Archaea phyla.
559 **b**, Summary of proteomics evidence of archaeal histones, including the presence of modifications,
560 tails, coverage, fraction of lysines identified, and isoelectric points. Human Histone H3 and H4 are
561 included for reference. The alignments at the bottom depict the position of lysine modifications in
562 the globular part of *Methanospirillum stamsii* and *Methanobrevibacter cuticularis* HMfB histones
563 (modified residues in bold). **c**, Archaeal HMfB histones with *N*-terminal tails (at least 10 aa before a
564 complete globular domain), sorted by frequency of lysine residues in the tail and color-coded ac-
565 cording to taxonomy (same as panel A). Amino-acid sequences shown for selected examples. The
566 dotted line indicates the median frequency of lysines in canonical eukaryotic H3 and H4 histone

567 tails. Source data available in Supplementary Table 2. **d**, Mass spectra of three modified archaeal
568 peptides, representing the relative abundance of fragments at various mass-to-charge ratios (m/z).
569 Spectra were annotated using IPSA. b and y ions and their losses of H₂O are marked in green and
570 purple, respectively; precursor ions are marked in dark grey. Unassigned peaks are marked in light
571 grey. Some labels have been omitted to facilitate readability.

572 **Figure 3. Taxonomic distribution of chromatin-associated gene classes.** **a**, Summary of the sev-
573 en classes of genes with chromatin-related activity covered in our survey: histone-specific hPTM
574 writers (acetylases and methyltransferases), erasers (deacetylases and demethylases), readers,
575 remodellers, and chaperones. **b**, Percentage of surveyed taxa containing homologs from each chro-
576 matin-associated gene class, for eukaryotes (top), archaea, bacteria, and viruses (bottom). Species-
577 level tables are available in Supplementary Fig. 3. **c**, Number of eukaryotic genes classified in each
578 of the chromatin-associated modification enzymes, readers, remodellers, and chaperones. **d**, Over-
579 lap between the taxon-level phylogenetic distribution of histones and chromatin-associated domains
580 in archaea and four bacterial phyla, measured using the Jaccard index. **e**, Number of genes encoding
581 writer, eraser, reader and remodeller domains, per species.

582 **Figure 4. Origin and evolution of chromatin-associated gene families.** **a**, Summary of phyloge-
583 netic affinities of the eukaryotic homologs of gene classes that are also present in prokaryotes. For
584 each gene family, we evaluate whether it is phylogenetically closer to a majority ($\geq 50\%$) of eukary-
585 otic sequences from a different orthogroup (indicating intra-eukaryotic diversification), or to se-
586 quences from Bacteria or Archaea. **b**, Left, gene tree of eukaryotic and prokaryotic Sirtuin
587 deacetylases, showcasing an example of a eukaryotic family that diversified within eukaryotes
588 (SIRT6) and another one with close relatives in Asgard archaea (SIRT7). Right, gene tree of KAT14
589 acetylase, a eukaryotic orthogroup with bacterial origins. Statistical supports (UF bootstrap) are
590 shown at selected internal nodes of the highlighted clades. **c**, Evolutionary reconstruction of hPTM
591 writer and eraser gene families, remodellers, and histone chaperones along the eukaryotic phyloge-
592 ny, including the number of genes present in the last eukaryotic common ancestor (LECA). Barplots
593 indicate the number of orthologs of each gene family present at the LECA (at 90% posterior proba-
594 bility; see Methods) and whether the presence of a given orthogroup at LECA is supported by its
595 conservation in various early-branching eukaryotic lineages (Amorphea, Discoba, Diaphoretickes
596 and others). The list of ancestral gene families below each plot is non-exhaustive. Two ancestral
597 gene counts are provided: all families at presence probability above 90%, and, in brackets, the sub-
598 set of these that is present in at least two of the main eukaryotic early-branching lineages
599 (Amorphea, Diaphoretickes, and Discoba). Source data in Supplementary Table 5. **d-e**, Recon-
600 structed evolutionary origins of the different subunits of the Polycomb repressive complexes (PRC2
601 and PRC1) and Trithorax-group complexes (KMT1 to 5). **f-h**, Side-by-side comparison of the pres-
602 ence of individual hPTM marks and various subunits of the Polycomb and Trithorax complexes, as
603 well as other hPTM writers, responsible for their deposition.

604 **Figure 5. Evolution of chromatin readers and capture of chromatin proteins by transposable**
605 **elements and viruses.** **a**, Evolutionary reconstruction of reader gene families along the eukaryotic
606 phylogeny, highlighting the number of gains along the eukaryotic phylogeny (at 90% posterior
607 probability). The Euler diagram at the top shows the overlap between presence of chromatin-
608 associated catalytic domains and readers. The barplot at the left indicates the number of orthologs of

609 each gene family present at the LECA and whether their presence is supported by its conservation in
610 various early-branching eukaryotic lineages (Amorphea, Discoba, Diaphoretickes, and others). Pie
611 plots at the right summarize the number of orthogroups from each gene family gained within select-
612 ed lineages: Metazoa, Holomycota, Viridiplantae and SAR+Haptophyta. **b**, Number of reader or
613 catalytic orthogroups gained at each node in the species tree, for selected species. Source data in
614 Supplementary Table 5. **c**, Networks of protein domain co-occurrence for Chromo and PHD read-
615 ers. Each node represents a protein domain that co-occurs with Chromo or PHD domains, and node
616 size denotes the number of co-occurrences with either Chromo or PHD. Edges represent co-
617 occurrences between domains. Groups of frequently co-occurring protein domains have been manu-
618 ally annotated and color-coded, which has revealed sub-sets of retrotransposon and DNA transpos-
619 on-associated domains. **d**, Number of chromatin-related eukaryotic genes fused with transposons
620 grouped by gene family (left), including the fraction that are classified as valid gene models based
621 on expression and assembly data (centre); and the number of species where each type of fusion is
622 found (right). The number of fusion events are colored according to their similarity with known
623 DNA transposons (red) or retrotransposons (orange) from the Dfam database (see Methods). (*) The
624 ‘Chromo’ category excludes genes containing other chromatin-associated protein domains such as
625 SNF2_N (listed separately as ‘Chromo+SNF2_N’, which includes remodellers with the domain of
626 unknown function DUF1087, which is also common in DNA transposons). **e**, Selected examples of
627 transposon fusion domains classified by orthogroup, including their archetypical protein domain
628 architecture, homology to transposon class, their phylogenetic distribution, and number of fusion
629 genes. Only orthogroups with at least one valid gene model are listed. Source data available in Sup-
630 plementary Table 6. **f**, Example tree of Chromo readers, highlighting genes with fused TE-
631 associated domains and their consensus domain architectures. **g**, Fraction of viral genomes contain-
632 ing homologs from each chromatin gene family, for nucleocytoplasmic giant DNA virus families
633 (top) and other taxa containing histone domains (Nudiviridae, Polydnviridae; bottom). **h**, Phyloge-
634 netic analysis of histone domains, with a focus on viral homologs. Statistical supports (approximate
635 Bayes posterior probabilities) are shown for the deepest node of each canonical eukaryotic or
636 archaeal histone clade. The inset table summarizes the presence of doublet histone genes per lineage.
637 **i**, Number of viral homologs in each chromatin-associated gene family, classified according to their
638 closest cellular homologs (eukaryotes, bacteria or archaea) in phylogenetic analyses (see Methods).
639 Source data available in Supplementary Table 6.

640 **Figure 6. Chromatin evolution and eukaryogenesis.** **a**, Summary of events in chromatin evolu-
641 tion prior to, during and after the origin of eukaryotes. **b**, Number of chromatin-related gene fami-
642 lies and hPTM marks inferred to have been present at the LECA. Ancestral gene counts are indicat-
643 ed at >90% probability. For gene counts, numbers within bars indicate the subset of families present
644 in at least two of the most deeply-sampled early-branching eukaryotic lineages (Amoropha,
645 Diaphoretickes, and Discoba). For hPTMs, the ancestral counts have been inferred using Dollo par-
646 simony assuming a Diaphoratickes – Amorphea split at the root of eukaryotes, and numbers within
647 bars indicate the number of hPTMs whose ancestral presence is supported by more than one species
648 at both sides of the root. **c**, hPTMs inferred to be present in the last eukaryotic common ancestor
649 (LECA) based on Dollo parsimony. Only amino-acid positions conserved in all eukaryotes in our
650 dataset are shown. Asterisks indicate modifications whose presence at the LECA is supported by

651 just one species at either side of the root. The inferred LECA presence of known writing/erasing
652 enzymes associated to these hPTM is indicated.

653

654 **SUPPLEMENTARY FIGURES LEGENDS**

655 **Supplementary Fig. 1. Histone classification and evolution. a**, Primary and secondary alignments
656 of histone-fold containing proteins classified as canonical H2A, H2B, H3 and H4, based on identity
657 to reference sequences in HistoneDB⁴⁸. Pie plots represent the number of alignments to HistoneDB-
658 annotated sequences, for the entire dataset (prokaryotic, eukaryotic and viral sequences, large pie
659 plots in the inset) and the eukaryotic subset (smaller plots in the inset). For those proteins that align
660 to more than one canonical histone or major variant (macroH2A, H2A.Z or cenH3), the scatter plots
661 represent the relative identity between the primary (horizontal axis) and secondary alignment(s)
662 (vertical axis). **b**, Aggregated counts of histone gene pairs, classified according to histone type and
663 orientation. **c**, Presence of histone variants (left) and number of collinear pairs of histone-encoding
664 genes (right) per species, classified according to their histone types and relative orientation (head-
665 to-head, hh; head-to-tail, ht; and tail-to-tail, tt). Source data available in Supplementary Table 2.
666 Histone variant classification is based on the highest-scoring HMM profile from HistoneDB. Aster-
667isks colors in the macroH2A column indicate species where histone-less Macro domains ortholo-
668gous to the macroH2A genes are found (see panel d). Lighter colors in the variant classification
669 indicate ambiguously classified histones (i.e. cases in which the highest-scoring HMM profile ex-
670 hibited a low bitscore, defined as a probability below 0.05 in the profile-wise distribution function
671 of scaled bitscores; or cases in which the first-to-second ratio between high scoring profiles was
672 below 1.01). **d**, Alignments of putatively conserved histone *N*-tails in archaea. Conserved amino-
673 acids are color-coded according to chemical properties. Dots next to species names are color-coded
674 according to taxonomy (same as **Fig. 2c**). **e**, Phylogenetic analysis of the Macro motif of macroH2A
675 histones across eukaryotes, highlighting the macroH2A ortholog group (green), and, within this
676 group, Macro-containing genes lacking histone domains (orange), and their protein domain archi-
677 tectures.

678 **Supplementary Fig. 2. Histone post-translational modifications. a**, Proteomics detection cover-
679 age (% of amino acids), number of hPTMs and number of hPTMs per covered position, for the best-
680 covered histone in each species in our proteomics survey. **b**, Number of samples in which each his-
681 tone-matching peptide with post-translational modifications (peptide spectral matches defined by
682 *Proteome Discoverer*) has been identified, per species. For each species, we report the percentage of
683 modified peptides found in more than one replicate. **c**, Number of samples in which histone-
684 matching modified peptide has been identified, across all the samples from this study. The tree pie
685 charts represent these distributions for all hPTMs, acetylations, and methylations. **d**, Evidence of
686 hPTM conservation in the major histone variants H2A.Z and macroH2A (conserved positions only),
687 as well as any position in the linker histones H1.

688 **Supplementary Fig. 3. Gene family counts. a-c**, Number of taxa within each lineage that contain
689 chromatin-associated genes, for archaeal, bacterial (per phyla) or viral (per family) genomes. Num-
690 bers indicate the exact number of taxa. **d**, Number of genes encoding core domains that define

691 chromatin-associated gene families per eukaryotic genome/transcriptome. Numbers indicate exact
692 number of proteins.

693 **Supplementary Fig. 4. Evolutionary reconstruction and domain architecture conservation. a,**
694 Species tree of eukaryotes used in the ancestral reconstruction analysis, with branch lengths cali-
695 brated to the gain/loss rates of Pfam domains (see Methods). Available in Supplementary Table 1. **b,**
696 Conservation of archetypical protein domain architectures across orthogroups, in acetylases,
697 deacetylases, methyltransferases, demethylases, remodellers and chaperones. In each heatmap, we
698 indicate the fraction of genes within an orthogroup (rows) that contain a specific protein domain
699 (columns). Domains in bold are catalytic (black) or reader (purple) functions. At the right of each
700 heatmap, we summarize the presence/absence profile of each orthogroup across eukaryotic lineages
701 (as listed in Fig. 1a).

702 **Supplementary Fig. 5. Evolution of the hPTM reader toolkit. a,** Pie plot representing the num-
703 ber of genes classified as part of the catalytic (acetylases, deacetylases, methyltransferases,
704 demethylases, remodellers or chaperones) or reader families, or as both. The barplot at the right
705 shows the most common reader domains in genes classified with both reader and catalytic func-
706 tions. **b,** Pie plot representing the number of reader domain-encoding genes classified according to
707 whether they contain one type of reader domain (e.g., PHD) or more than one (e.g., PHD +
708 PWWP). The barplot at the right shows the most common combinations of reader domains among
709 genes with multiple reader domains. **c,** Summary of gene family gains per reader family, with ex-
710 ample cases highlighted in selected nodes. Node size is proportional to number of gains at 90%
711 probability.

712 **Supplementary Fig. 6. Transposon-chromatin gene fusions. a,** Number of candidate fusion genes
713 classified by the level of gene model validation evidence, based on contiguity of the gene model
714 over the genome assembly (i.e. lack of poly-N stretches in the genomic region between the TE- and
715 chromatin-associated domains), evidence of expression, and evidence of contiguous expression (see
716 inset at the right). **b,** Summary of candidate gene fusions within each chromatin-associated gene
717 family, divided by gene family. For each gene, we indicate their similarity to known TE families,
718 presence of TE-associated domains, the evidence of gene model validity, and information on their
719 gene structure (whether they are monoexonic or are located in clusters with other fusion genes).
720 Source data available in Supplementary Table 6. **c,** Number of species with at least one valid fusion,
721 divided by gene family. **d,** Mapping positions of RNA-seq reads supporting candidate gene-
722 transposon fusions (selected examples from Fig. 5e). For each fusion, we show reads spanning the
723 region along the spliced transcript that fully covers the transposon-associated domains (highlighted
724 in green), the chromatin-associated domains, and the inter-domain region. Uninterrupted stretches
725 of mapped positions between domains indicate the validity of a domain co-occurrence. For clarity
726 purposes, reads mapping entirely within a single domain have been excluded from this visualiza-
727 tion.

728 **Supplementary Fig. 7. Chromatin proteins in viruses. a-c,** Selected gene trees highlighting ex-
729 amples of eukaryotic- and prokaryotic-like viral homologs. **d,** Number of viral genes of each chro-
730 matin-associated gene family, classified according to their closest neighbours from cellular clades in
731 gene tree analyses based on phylogenetic affinity scores (see Methods). Within each gene family,

732 viral sequences are classified according to their PFAM domain architecture – the most common
733 architecture being single-domain in most gene families except for remodellers and BIR readers. **e**,
734 *Id.*, but classifying viral genes according to their phylogenetic affinity to eukaryotic orthology
735 groups. Source data available in Supplementary Table 6.

736 **Supplementary Material 8. Phylogenetic analyses.** Collection of gene trees used to identify
737 orthology groups for the eukaryotic chromatin toolkit. UFBS bootstrap supports rare indicated at
738 each node. An annotated eukaryotic species tree is also included.

739 **Supplementary Material 9. Peptide sequences.** Collection of peptide sequences used to build
740 gene trees of the eukaryotic chromatin toolkit.

741

742

743

744 **SUPPLEMENTARY TABLES LEGENDS**

745 **Supplementary Table 1. Taxon sampling. a**, List of eukaryotic species used in the comparative
746 genomic analyses, including species abbreviations, data sources for genome or transcriptome as-
747 semblies and annotations, and their taxonomic classification. **b**, List of gene expression datasets
748 (SRA accession numbers) used for gene model validation analyses of candidate fusion genes. **c**, List
749 of histone post-translational modification proteomics datasets used in this study (PRIDE accession
750 numbers).

751 **Supplementary Table 2. Histone clusters and classification. a**, Pairs of collinear histone-
752 encoding genes, including their genomic coordinates and relative orientation. **b**, List and sequences
753 of archaeal HMfB histones with N-terminal tails (at least 10 aa before a complete globular domain).
754 **c**, Classification of histone variants across eukaryotes.

755 **Supplementary Table 3. hPTM conservation. a-g**, Table of hPTMs identified in histones of the 26
756 eukaryotic species used in the comparative proteomics analysis, separated by histone type (canoni-
757 cal and major variants: H2A, H2B, H3, H4, macroH2A, H2A.Z, and H1). Each entry corresponds to
758 a modified peptide, for which we specify modification coordinates along the peptide and relative to
759 the consensus histone sequence (if available). We also indicate whether each peptide can be unique-
760 ly mapped to a conserved or non-conserved region in a canonical histone, or to specific histone var-
761 iants. These tables also include entries for hPTMs reported in the literature (indicated as a cited
762 source or as a specific UNIPROT entry; see Methods for a list of sources); in these cases, source
763 peptides and associated data may not be available. **h**, hPTMs in Archaea.

764 **Supplementary Table 4. Gene family analysis. a**, List of gene classes analyzed in the comparative
765 genomics analyses, including the PFAM protein domains used to retrieve homologs and search pa-
766 rameters. **b**, List of transposon-associated PFAM domains surveyed in the analyses of transposon-
767 chromatin gene fusions.

768 **Supplementary Table 5. Evolution of the chromatin machinery in eukaryotes. a**, Summary of
769 gene family evolutionary patterns in eukaryotes ($n = 1,713$ orthogroups). For each orthogroup, we
770 indicate its gene and functional class, the number of members, species where it is present, and ma-
771 jor eukaryotic lineages (Amoebozoa, Opisthokonta+Breviatea+Apusozoa, CRuMs,
772 Ancyromonadida, Mala-wimonadidae, Archaeplastida+Cryptista, SAR+Haptista,
773 Hemimastigophora, Discoba, and Metamonada), the probability of presence at the last eukaryotic
774 common ancestor, the phylogenetic affinity of their closest homologs (other eukaryotic orthogroups,
775 bacteria, archaea or viruses) and their average frequency amongst the 10 nearest neighbours of its
776 member gene in phylogenetic trees ('Phylogenetic affinity score', see Methods); as well as its con-
777 sensus protein domain architecture (present in at least 25% of its members). We also indicate the
778 gene symbols of members from four model species: *H. sapiens*, *D. melanogaster*, *S. cerevisiae*, and
779 *A. thaliana*. **b-c**, Probability of gain and loss of each gene family at extant and ancestral nodes along
780 the eukaryotic phylogeny. **d**, Orthogroup assignments per gene.

781 **Supplementary Table 6. Transposon fusions and viral homology. a**, List of candidate fusions
782 between chromatin-associated genes and transposons, including the phylogenetic classification of
783 each gene (orthogroup), protein domain architectures, and the transcriptomics-level and gene mod-
784 el-level evidence supporting each fusion. **b**, List of chromatin-associated genes encoded by viral
785 genomes, including their species of origin and a summary of their phylogenetic embedding among
786 cellular species (specifically, which are its closest homologs in cellular genomes and the fraction of
787 phylogenetic nearest neighbours they represent, the closest eukaryotic gene family among those
788 close to eukaryotic genes in the gene trees, and the distance to the closest cellular homolog).

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817 **Methods**

818 *Eukaryotic cell culture and tissue sources*

819 *Capsaspora owczarzaki* strain ATCC30864 filopodial cells were grown axenically in 5 ml flasks
820 with ATCC medium 1034 (modified PYNFH medium) in an incubator at 23°C (Sebé-Pedrós et al.,
821 2013a).

822 *Corallochytrium limacisporum* strain India was axenically grown in Difco Marine Broth medium at
823 23°C, *Creolimax fragrantissima* strain CH2 was axenically grown in Difco Marine Broth medium at
824 12°C, *Spizellomyces punctatus* strain DAOM BR117 was axenically grown in (0,5% yeast extract,
825 3% glycerol, 1g/L K₂HPO₄, 0,5% EtOH) medium at 17°C, *Thecamonas trahens* strain ATCC50062
826 was grown in ATCC medium: 1525 Seawater 802 medium, *Chlamydomonas reinhardtii* strain CC-
827 503 cw92 mt+ was axenically grown in Gibco TAP medium at 29°C, *Guillardia theta* strain
828 CCMP2712 was axenically grown in L1+500uM NH₄Cl medium at 18°C, *Emiliana huxleyi* strain
829 CCMP1516 was grown in L1-Si medium at 18°C, *Thalassiosira pseudonana* strain CCMP1335 was
830 axenically grown in L1 medium at 18°C, *Bigelowiella natans* strain CCMP2755 was axenically
831 grown in L1-Si medium at 23°C, *Naegleria gruberi* strain ATCC30224 was axenically grown in
832 ATCC medium 1034 (modified PYNFH medium) at 29°C, *Gefionella okellyi* strain 249 was grown
833 in 15% Water Complete Cereal Grass Media (WC□CGM3) at 18°C and *Fabomonas tropica* strain
834 NYK3C was grown in L1 + YT medium at 18°C. All cells were grown in 250 ml culture flasks.

835 In addition, we used frozen tissues/cells from the following species: *Homo sapiens* (ES cells, cour-
836 tesy of Cecilia Ballaré, CRG), *Physcomitrella patens* (strain Gransden 2004, vegetative stage, cour-
837 tesy of Josep Casacuberta, CRAG-CSIC), *Sycon ciliatum* (adult sponges sampled from Bergen,
838 Norway, courtesy of Maja Adamska, ANU) and *Phytophthora infestans* (strain T30-4, courtesy of
839 Harold J.G.Meijer, Wageningen University).

840

841 **Archaeal cell culture**

842 Cultures of *Methanobrevibacter cuticularis* DSM 11139, *Methanospirillum stamsii* DSM 26304 and
843 *Methanosarcina spelaiei* DSM 26047 were purchased from the Deutsche Stammsammlung von
844 Mikroorganismen und Zellkulturen GmbH (DSMZ), Braunschweig, Germany. Cultures were grown
845 in closed batch in 50mL of defined media in 120mL serum bottles (La-Pha-Pack, Langerwehe,
846 Germany). Growth was monitored as OD (600 nm; Analytik Jena, Specord 200 plus).
847 *Methanobrevibacter cuticularis* was grown in modified *Methanobrevibacter cuticularis* medium
848 DSMZ 734a (DSMZ 2014) omitting bovine rumen fluid, yeast extract and Na-resazurin at 1.5 bar
849 overpressure H₂CO₂ (20 vol.-% CO₂ in H₂) at 37°C. As soon as a change in OD was observed, a
850 constant agitation at 90rpm was applied. *Methanospirillum stamsii* was grown in modified
851 *Methanobacterium* medium DSMZ 119 (DSMZ 2017) omitting sludge fluid, yeast extract and Na-
852 resazurin at 1 bar overpressure H₂CO₂ (20 vol.-% CO₂ in H₂) at 29°C, under constant agitation at
853 90rpm. *Methanosarcina spelaiei* was grown in modified *Methanosarcina barkeri* medium DSMZ
854 120a (DSMZ 2014) omitting yeast extract and Na-resazurin at 1.5 bar overpressure H₂CO₂ (20 vol.-
855 % CO₂ in H₂) at 33°C, under constant agitation at 90rpm. All gases were obtained from Air Liquide
856 GmbH, Schwechat, Austria. *Nitrososphaera viennensis* EN76 was grown in continuous culture in a
857 bioreactor as previously described¹⁰⁰.
858 Cells were harvested via centrifugation at 21,000xg 4°C 1h (Thermo scientific, Sorvall Lynx 4000
859 centrifuge), the supernatant discarded and the resulting pellet resuspended in 1ml of spent medium,
860 followed by another round of centrifugation at 21,000xg 4°C for 1h (Eppendorf, Centrifuge 5424R).
861 Pellets were stored at -70°C. All archaeal histones were extracted as described below.

862

863 **Histone acid extraction**

864 Starting material was a pellet of 50-100M cells (washed once with cold PBS) or a flash-frozen tis-
865 sue homogenate in liquid nitrogen using a ceramic mortar grinder. Cells were washed first in 10ml
866 of buffer I (10 mM TrisHCl pH 8, 10 mM MgCl₂, 0.4M Sucrose). After 5min incubation, samples
867 were centrifuged at 8.000g for 20min at 4°C and supernatant was removed. The resulting pellet was
868 resuspended in 1.5ml of Buffer II (10 mM TrisHCl pH 8, 10 mM MgCl₂, 0.25M Sucrose, 1% Triton
869 X-100, 1% Igepal Ca-630) and incubated 15min on ice. In specific cases, cells at this stage were
870 broken using a 2ml Dounce homogenizer (with Pestle B) or with a 20G syringe. Then samples were
871 centrifuged at 15.000g for 10min at 4°C and supernatant was removed. The resulting pellet was then
872 slowly resuspended in 300μL of Buffer III (10 mM TrisHCl pH 8, 2 mM MgCl₂, 1.7M Sucrose, 1%
873 Triton X-100) and then resulting resuspended nuclei were layered on top of another 300μL of Buff-
874 er III. Sample was centrifuged at 20.000g for 1h at 4°C and supernatant was removed, resulting in a
875 nuclear pellet ready for acid histone extraction. All buffers were supplemented with spermidine
876 (1:1000), beta-mercaptoethanol (1:1000), protease inhibitors (1x cOmplete cocktail Roche
877 #11697498001, 1mM PMSF, 1:2000 Pepstatin), phosphatase inhibitors (1x phoSTOP cocktail
878 Roche #4906845001) and deacetylase inhibitors (10mM Sodium butyrate).
879 For samples processed using a high-salt + HCl extraction protocol^{101,102}, the pellet was resuspended
880 in 500μL of High Salt Extraction Buffer (20 mM TrisHCl pH 7.4, CaCl₂ 1M and protease, phosphatase
881 and deacetylase inhibitors, same as above). Sample was incubated on ice for 30min and then

882 pure HCl has added to a final 0.3N concentration (12.82 μ L to the initial 500 μ L). Samples were in-
883 cubated for at least 2h on a rotor at 4°C and then centrifuged at 16.000g for 10min at 4°C to remove
884 cellular/nuclear debris. The resulting supernatant containing solubilized histones was transferred to
885 a clean 1.5ml tube and Trichloroacetic Acid (TCA) was added drop-wise to 25% final concentration
886 (171 μ L TCA to an approximate initial 513 μ L sample) and left overnight at 4°C to precipitate his-
887 tones. Samples were then centrifuged at 20.000g for 30min at 4°C and the supernatant removed. The
888 pellet was then washed twice with 500 μ L of cold acetone and then dried for 20min at room temper-
889 ature. Finally, clean histone pellets were resuspended in 30-50 μ L of ultrapure water. Protein con-
890 centration in the sample was measured using BCA and extraction was examined using an SDS-
891 PAGE protein gel with Coomassie staining.

892 For samples processed using H₂SO₄¹⁰², the protocol was exactly the same except that 400 μ L 0.4N
893 H₂SO₄ (freshly diluted) was used instead, with a similar incubation time of at least 2h at 4°C.

894

895 ***Histone chemical derivatization***

896 Histones samples were quantified by the BCA method and 10 μ g of each sample were derivatized
897 with propionic anhydride, digested with trypsin and derivatized again with phenylisocyanate as pre-
898 viously described⁴⁹. Briefly, samples were dissolved in 9 μ L of H₂O and 1 μ L of triethyl ammoni-
899 um bicarbonate was added to bring the pH to 8.5. The propionic anhydride was prepared by adding
900 1 μ L of propionic anhydride to 99 μ L of H₂O and 1 μ L of propionic anhydride solution was added
901 immediately to the samples with vortexing and incubation for 2 minutes. The reaction was
902 quenched with 1 μ L of 80mM hydroxylamine and samples were incubated at room temperature for
903 20 minutes. Tryptic digestion was performed for 3 h with 0.1 μ g trypsin (Promega Sequencing
904 Grade; Madison, WI) per sample. A 1% v/v solution of phenyl isocyanate (PIC) in acetonitrile was
905 freshly prepared and 3 μ l added to each sample (17 mM final concentration) and incubated for 60
906 min at 37 °C. Samples were acidified by adding 50 μ L of 5% formic acid, vacuum dried and desalt-
907 ed with C18 ultramicrospin columns (The Nest Group, Inc, Southborough, MA).

908

909 ***Liquid Chromatography-Tandem Mass Spectrometry Sample Acquisition***

910 A 2- μ g aliquot of the peptide mixture was analyzed using a LTQ-Orbitrap Fusion Lumos mass spec-
911 trometer (Thermo Fisher Scientific, San Jose, CA) coupled to an EASY-nLC 1000 (Thermo Fisher
912 Scientific, San Jose, CA) with both collision induced dissociation (CID) and high energy collision
913 dissociation (HCD) fragmentation.

914 Peptides were loaded directly onto the analytical column and were separated by reversed-phase
915 chromatography using a 50-cm column with an inner diameter of 75 μ m, packed with 2 μ m C18
916 particles spectrometer (Thermo Scientific, San Jose, CA, USA) with a 90 min chromatographic gra-
917 dient. The mass spectrometer was operated in positive ionization mode using a data dependent ac-
918 quisition method. The “Top Speed” acquisition algorithm determined the number of selected pre-
919 cursor ions for fragmentation.

920

921 ***Mass-spectrometry Data Analysis***

922 Acquired data were analyzed using the Proteome Discoverer software suite (v2.0, Thermo Fisher
923 Scientific), and the Mascot search engine (v2.6, Matrix Science¹⁰³) was used for peptide identifica-
924 tion using a double-search strategy. First, data were searched against each organism protein data-
925 base plus the most common contaminants considering Propionylation on *N*-terminal, Propionylation
926 on Lysines and Phenylisocyanate on *N*-terminal as variable modifications. Then a new database was
927 generated with the proteins identified in the first search,, and a second search was done considering
928 Propionylation on *N*-terminal, Propionylation on Lysines, Phenylisocyanate on *N*-terminal, Dime-
929 thyl lysine, trimethyl lysine, propionyl + methyl lysine, acetyl lysine, crotonyl lysine as variable
930 modifications. Precursor ion mass tolerance of 7 ppm at the MS1 level was used, and up to 5 missed
931 cleavages for trypsin were allowed. False discovery rate (FDR) in peptide identification was set to a
932 maximum of 5%. The identified peptides were filtered by mascot ion score higher than 20 and only
933 PTMs with a localization score ptmRS¹⁰⁴ higher than 45 were considered. The raw proteomics data
934 have been deposited to the PRIDE¹⁰⁵ repository with the dataset identifier PXD031991.

935

936 *Analysis of hPTM conservation*

937 *Identification of canonical and variant histones.* We classified histone protein domains from a data-
938 base of eukaryotic, prokaryotic and viral sequences (see details below) according to their similarity
939 to known canonical (H2A, H2B, H3, H4) and variant histones (e.g., H2A.Z, macroH2A, cenH3 or
940 H3.3), as well as other gene families with histone-like protein folds (e.g., the transcription factors
941 DR1, DRAP1, NFYB/C, POLE3/4, SOS, TAF, or CHRAC). To that end, we used *diamond* to per-
942 form local alignments of each histone domain against (i) a set of curated histone variants obtained
943 from HistoneDB 2.0⁴⁸, and (ii) annotated each domain according to the best hit in the reference da-
944 tabase, which allowed us to classify histone fold-containing proteins as canonical histones (H2A,
945 H2B, H3, H4) or their main variants (H2A.Z, macroH2A and cenH3). This best-hit strategy per-
946 forms well in distinguishing canonical histones from each other, as well as each canonical histone
947 from its main variants (H3 from cenH3, and H2A from H2A.Z and macroH2A; **Supplementary**
948 **Fig. 1a**).

949 Then, we built a graph of pairwise similarity between histones, with edges weighted by the align-
950 ment bitscore (discarding edges with bitscore < 20). We created visualisations of each connected
951 component in this graph using the spring layout algorithm implemented in the *networkx* 2.4 Python
952 library (100 iterations, weighted by alignment bitscore)¹⁰⁶. We selected the four connected compo-
953 nents in the graph that matched the four canonical eukaryotic histones (H2A, H2B, H3, H4; discard-
954 ing edges with bitscore < 20), retrieved the protein sequences for each of them, aligned them using
955 *mafft* (E-INS-i mode, 1,000 iterations)¹⁰⁷, and built phylogenetic trees with *IQ-TREE* 2.1.0 (*-fast*
956 *mode*)¹⁰⁸.

957 *Identification of hPTM homology.* We retrieved the protein sequences of the canonical histones
958 identified in each of the 26 species and we used them for the proteomic analysis of hPTMs, and
959 aligned them using *mafft* (*G-INS-i* mode, up to 10,000 refinement iterations). For this subset of spe-
960 cies, histone class identity was cross-referenced with the HistoneDB search tool. Then, we manually
961 aligned the peptides mapping onto these proteins to identify the position of each hPTM along a con-
962 sensus alignment. In the case of H3, H4, and macroH2A, the majority of alignment positions were
963 conserved across most eukaryotes in our dataset, and we used a consensus numbering scheme. In

964 the case of H2A, H2A.Z, and H2B, non-conserved insertions and deletions at the N-terminal tail
965 precluded the use of a paneukaryotic numbering scheme. Instead, we reported hPTM positions
966 based on the human homolog (if possible), or relative to taxonomically restricted conserved posi-
967 tions. In cases where position-wise homology could not be established, we grouped multiple amino-
968 acids into stretches of unclear homology, which we report separately from conserved positions
969 (question mark symbols in **Fig. 1**). The complete list of hPTMs and their position-wise coordinates
970 relative to the consensus alignment is available in **Supplementary Table 3**.

971 Furthermore, we also reported the presence (in any position) of modifications in less-conserved
972 histone variants, as well as the linker histone H1.

973 In addition to the 19 used in our proteomics survey, we also included previously published hPTM
974 data from the following species (**Supplementary Table 1c**): the brown alga *Ectocarpus*
975 *siliculosus*⁴³, the diatom *Phaeodactylum tricornutum*¹⁰⁹, the ciliate *Tetrahymena thermophila*^{46,110-}
976 ¹¹², the ascomycete *Neurospora crassa*¹¹³, *Saccharomyces cerevisiae* and *Schizosaccharomyces*
977 *pombe*⁴⁶, and the plant *Arabidopsis thaliana*¹¹⁴⁻¹¹⁶. When available in public repositories, we re-
978 analysed these datasets using the strategy described above. Finally, we also complemented our own
979 proteomics data using previously published hPTM data from *Homo sapiens*^{46,117-120} and
980 *Capsaspora owczarzaki*⁴².

981 **Comparative genomics analysis of chromatin-associated proteins**

982 *Data retrieval.* We identified homologs of gene families associated with eukaryotic chromatin, us-
983 ing a database of predicted proteomes from a selection of eukaryotic species from all major super-
984 groups ($n = 172$ species; see **Supplementary Table 1** for their taxonomic classification and data
985 sources), as well as archaeal and viral peptides available in the NCBI *nr* peptide collection (as of
986 25th of April, 2020) and bacterial peptides available in RefSeq (release 99, 11th May, 2020). The
987 database of viral sequences was complemented with peptides from 501 genomes of
988 nucleocytoplasmic large DNA viruses¹²¹.

989 *Gene family searches.* We defined 61 gene classes associated with eukaryotic chromatin, based on
990 HMM models obtained from the Pfam database (release 33.0)¹²². This list included canonical and
991 linker histones ($n = 2$ families), chromatin-specific lysine acetylases ($n = 5$), deacetylases ($n = 2$),
992 methyltransferases ($n = 2$), demethylases ($n = 2$), chromatin readers ($n = 16$), remodellers ($n = 1$)
993 and chaperones ($n = 13$), as well as multiple families associated with the Polycomb complexes ($n =$
994 18). The complete list of gene families, including the associated HMM models, is available in **Sup-**
995 **plementary Table 4**.

996 For each gene family, we retrieved all homologs from the eukaryotic, archaeal, bacterial and viral
997 databases using the *hmmsearch* tool from the *HMMER* 3.3 toolkit¹²³ and the gathering threshold
998 defined in each Pfam HMM model. We recorded the taxonomic profile of each homolog.

999 *Orthology identification.* We aimed to identify groups of orthologs within each of the 61 chromatin-
1000 associated gene families using targeted phylogenetic analyses. We followed the following strategy
1001 for each of the 59 sets of eukaryotic genes. First, we partitioned each set into one or more homology
1002 groups based on pairwise local sequence alignments using *diamond* 0.9.36.137 (high sensitivity all-
1003 to-all search)¹²⁴, followed by clustering of the resulting pairwise alignments graph with *MCL*
1004 14.137 (--*abc* mode)¹²⁵, using low inflation values (see **Supplementary Table 4**) to favour inclu-

1005 sive groupings. Second, we performed multiple sequence alignments of each homology group with
1006 *mafft* 7.471¹⁰⁷ under the E-INS-i mode (optimised for multiple conserved regions), running up to
1007 10,000 refinement iterations. Third, we trimmed the resulting multiple sequence alignments using
1008 *clip-kit* 0.1 (*kplic-gappy* mode)¹²⁶. Fourth, we built phylogenetic trees for each trimmed alignment
1009 using *IQ-TREE* 2.1.0¹⁰⁸, selecting the best-fitting evolutionary model using its *ModelTest* module
1010 (according to the Bayesian Information Criterion) and using 1,000 UFBS bootstrap supports¹²⁷.
1011 Each tree was run for up to 10,000 iterations until convergence was attained (at the 0.999 correla-
1012 tion coefficient threshold, and for at least 200 iterations).

1013 Then, we parsed the species composition of each gene tree in order to identify groups of ortholo-
1014 gous proteins using the *POSSVM* pipeline¹²⁸. Specifically, we used the species overlap algorithm¹²⁹
1015 implemented in the *ETE* toolkit 3.1.1¹³⁰, which identifies pairs of orthologous genes in a phyloge-
1016 netic tree by examining the species composition of each subtree, and classifying internal nodes as
1017 paralogy nodes (if there is overlap in the species composition between each of its two descendant
1018 subtrees) or orthology nodes (if there is no overlap). Pairs of genes linked by an orthology node are
1019 then recorded as orthology pairs. In our analysis, we used an overlap threshold=0 (i.e. any species
1020 composition overlap between the two descendant subtrees is classified as a paralogy event). The
1021 resulting list of pairwise orthology relationships between genes was clustered into groups of
1022 orthologs (orthogroups) using *MCL*. We further annotated each orthogroup with a string denoting
1023 the gene symbols of the human proteins therein (if any).

1024 Overall, we classified 51,426 proteins from 61 gene classes (defined by protein structural domains),
1025 divided into 242 gene trees and 1,713 gene families (orthogroups). The source peptide sequences
1026 and gene trees used for these analyses are available in **Supplementary Material 7 and 8**.

1027 *Ancestral reconstruction of gene content.* We inferred the presence, gain and loss of each
1028 orthogroup along the eukaryotic tree of life, using a phylogenetic birth-and-death model¹³¹ imple-
1029 mented in *Count*¹³². This tool takes a numeric profile of gene family presence/absence in extant
1030 species (172 in our dataset) and a phylogenetic tree defining their evolutionary relationships, and
1031 infers the probabilities of gain and loss of each family at each ancestral node along the tree.

1032 First we trained the probabilistic model in *Count*. As a training set, we used a random sample of
1033 1,000 PFAM domains annotated in the 172 species of interest (restricting the sampling to domains
1034 present in at least 5% of species). The final model consists of gain, loss and transfer rates with two
1035 Γ categories each, and a constant duplication rate (given that we only recorded gene pres-
1036 ence/absence, duplication events are not included in our downstream analyses). This model was
1037 obtained in three sequential rounds of training, so as to sequentially add zero, one and two Γ cate-
1038 gories to each evolutionary rate. Each round consisted of up to 100 iterations, and stopped when the
1039 relative change in the model log-likelihood fell by 0.1% in two consecutive rounds. The final evolu-
1040 tionary rates and the Newick-formatted species tree used in this step are available in the **Supple-**
1041 **mentary Table 1 and Supplementary Fig. 3a**.

1042 Second, we calculated the posterior probability of gain, loss and presence of each orthogroup in our
1043 dataset with *Count*. The aggregated counts of gains and losses of the various classes of chromatin-
1044 associated proteins (acetylases, deacetylases, methyltransferases, demethylases, readers and
1045 remodellers) along the eukaryotic tree were obtained by summing the probabilities of gain, presence

1046 or loss of all orthogroups of a given class at each ancestral node. To investigate the evolutionary
1047 histories of specific orthogroups at a given node in the tree, we applied a probability threshold of
1048 0.9 (for presence) or 0.5 (to identify the most probable gain and loss node). The *Count* model was
1049 not able to calculate ancestral probabilities for a few orthogroups with widespread phylogenetic
1050 distributions, due to violations of the birth-and-death model (25 out of 1,713 families). In order to
1051 be able to report presence probabilities in the LECA for these orthogroups, we inferred their pres-
1052 ence in this ancestor using the Wagner parsimony procedure implemented in *Count* with a gain-to-
1053 loss penalty $g = 5$, and recorded their presence as binary values (0/1) accordingly.

1054 *Protein domain architecture analyses.* We annotated the Pfam domains present in each protein from
1055 the gene classes listed in **Supplementary Table 4**, using *Pfamscan* 1.6-3 and the Pfam 33.0 data-
1056 base¹²². We visualized the networks of protein domain co-occurrence from the point of view of the
1057 core domain(s) that define each gene class, using the *networkx* Python library (version 2.4)¹⁰⁶. Spe-
1058 cifically, we built a graph where each node represented ‘accessory’ domains (i.e. domains that co-
1059 occur with the ‘core’ domain that defines given gene class), node size reflected number of co-
1060 occurrences with the ‘core’ domain, and edges reflected co-occurrences between accessory do-
1061 mains. We identified communities of frequently co-occurring accessory domains using the label
1062 propagation algorithm implemented in *networkx* (*communities* submodule), which we used as a
1063 basis to manually annotate groups of co-occurring domains of interest (**Fig. 5C**). Network visualiza-
1064 tions were created using the *NEATO* spring layout algorithm from the *Graphviz* 2.40.1 Python li-
1065 brary¹³³.

1066 In parallel, we also recorded the presence of Pfam domains within individual orthogroups, and their
1067 taxonomic distribution.

1068 *Prokaryotic roots of the eukaryotic chromatin machinery.* We retrieved all eukaryotic domains from
1069 gene class shared with prokaryotes (Histones, Acetyltransf_1, GNAT_acetyltr_2, MOZ_SAS,
1070 Hist_deacetyl, SIR2, DOT1, SET, CupinJmjC, ING, MBT, PWWP and SNF2_N), collapsing identi-
1071 cal sequences at 100% similarity with *CD-HIT* 4.8.1¹³⁴, and identified their closest homologs
1072 amongst the corresponding archaea and bacteria protein domain sets, using *diamond* local align-
1073 ments (high sensitivity search). The archaeal and bacterial protein sets were also reduced with *CD-*
1074 *HIT* (at 95% and 90% sequence similarity, respectively). Each set of sequences was then partitioned
1075 into low-granularity homology clusters using the *MCL*-based strategy described above (inflation $I =$
1076 1.2), and a phylogenetic tree was then constructed from each homology cluster with *IQ-TREE* (as
1077 described above).

1078 Then, we mapped each eukaryotic gene to its orthogroup (obtained from eukaryotic-only analyses,
1079 see above) and used the distribution of phylogenetic distances from the prokaryotic+eukaryotic
1080 gene trees to classify them according to their similarity to (i) eukaryotic genes in other orthogroups,
1081 (ii) archaeal homologs, or (iii) bacterial homologs. Specifically, we used a majority-voting proce-
1082 dure in which we recorded the number of sequences of eukaryotic, archaeal or bacterial origin
1083 amongst the ten nearest neighbors of each gene (measuring intergenic distances as substitutions per
1084 site), and assigned the most common taxonomic group as the ‘closest’ homolog of that gene (mini-
1085 mum 50% agreement). This fraction is termed ‘Phylogenetic affinity score’ and reported in **Sup-**

1086 **plementary Table 5.** The pairwise distances were obtained from each gene tree using the
1087 cophenetic distance method in the *cophenetic.phylo* utility of the *ape* 5.4 *R* library¹³⁵.

1088 *Characterisation of fusions with transposon-associated domains.* We retrieved all classified genes
1089 from our eukaryotic dataset that contained transposon-associated Pfam domains (version 33.0), us-
1090 ing a list compiled from^{68,136} (complete list in **Supplementary Table 4**), totaling 823 candidate fu-
1091 sions from 91 species (listed in **Supplementary Table 6**). We annotated these genes to their most
1092 similar known TE element by aligning them against the Dfam 3.3 database¹³⁷ using the *tblastn* pro-
1093 gram in *BLAST* 2.2.31¹³⁸.

1094 We validated each candidate fusion using the following criteria: (i) contiguity of the gene model on
1095 the genome assembly, i.e., recording which genes were interrupted by poly-*N* stretches (which
1096 might indicate an incorrect gene model); (ii) evidence of expression in at least one sample from a
1097 range of publicly available transcriptomic experiments (from the NCBI SRA repository); (iii) evi-
1098 dence of contiguous expression, i.e., whether an expressed transcript had mapped reads along the
1099 entire region located between the ‘core’ and ‘TE-associated’ domains; (iv) we also recorded the
1100 number of exons per gene; and (v) located near any other candidate fusion gene in the genome.

1101 The list of SRA experiments used for these validation steps is available in **Supplementary Table 1**.
1102 This list includes 64 out of 91 species for which transcriptomics datasets are publicly available, and
1103 covers 768 out of the 822 TE fusion candidates (93%). RNA-seq read mapping was performed with
1104 *bwa mem* 0.7.17-r1188¹³⁹ using the complete set of spliced transcripts of each species as the refer-
1105 ence database. We used *bedtools* 2.29.2¹⁴⁰ to identify poly-*N* stretches in the genome assembly (as-
1106 sembly contiguity criterion). We identified regions of low coverage along the transcript sequence
1107 (expression contiguity criterion) using the *bedtools genomecov* utility, requiring that the coverage
1108 along both domains involved in each fusion and their intermediate regions be higher or equal to two
1109 reads.

1110 *Analysis of viral homologs.* We investigated the homology of the viral chromatin-associated genes
1111 (which included 19 out of 61 families present in our survey) using joint phylogenetic analyses of
1112 protein domains from virus, prokaryotic and eukaryotic genes. We used the same method described
1113 above to investigate the prokaryotic roots of eukaryotic gene classes: we aligned viral domains
1114 against a database of cellular homologs (high sensitivity *diamond* search), followed by low-
1115 granularity *MCL* clustering (inflation $I = 1.2$) and phylogenetic tree building (*IQ-TREE*). Then, we
1116 used the same majority-voting procedure described above to classify viral homologs according to
1117 their similarity to eukaryotic, archaeal or bacterial gene families based on their distribution of phy-
1118 logenetic distances. For viral genes that were most similar to eukaryotic genes, we used the same
1119 procedure to map them to their closest eukaryotic orthogroup.

1120 The complete list of viral genes and their phylogenetic annotation is available in **Supplementary**
1121 **Table 6**. Out of 2,163 viral genes in our dataset, 2,144 could be annotated as similar to a particular
1122 cellular group using this procedure (99.1%), and the majority of these genes had a high agreement
1123 in the annotations of their nearest neighbors (2,096 with $\geq 50\%$ agreement; 1,449 with $\geq 90\%$ agree-
1124 ment).

1125 In the case of viral histones, we built a separate phylogeny with a few modifications in our protocol:
1126 (i) we used additional viral genes obtained from⁷¹ as a reference; (ii) we omitted the *CD-HIT* reduc-

1127 tion and *MCL* partitioning steps, and jointly analyzed the entire set of homologs instead; and (iii) in
1128 the phylogenetic reconstruction step, we used the approximate Bayes posterior probabilities¹⁴¹ im-
1129 plemented in *IQ-TREE*.

1130 *Identification of archaeal N-terminal histone tails.* We retrieved all archaeal histone domains classi-
1131 fied belonging to the HMfB-like connected component in **Fig. 1b**, and retained those that fulfilled
1132 the following criteria: (i) contained a complete Cbfd_NFYB_HMF domain according to the
1133 *hmmscan* search (defined as an alignment starting at least at the 10th position of the HMM model,
1134 and up to the 55th position; the HMM model contains 65 positions); and (ii) the predicted tail (*N*-
1135 terminal to the core domain boundaries defined by *hmmscan*) was at least 10 residues long. 84
1136 genes passed these filters, including three *N*-terminal containing histones previously identified by
1137 Henneman *et al.*⁵⁵. A complete list is available in **Supplementary Table 2**. We manually examined
1138 the sequences of archaeal tails and aligned four sets of similar histones with *mafft G-INS-i* (**Sup-**
1139 **plementary Fig. 1d**). Alignments were plotted using the *msa* 1.24.0 library in *R*¹⁴².

1140

1141 *Data and Code Availability*

1142 The mass spectrometry proteomics data have been deposited to the ProteomeXchange Consortium
1143 via the PRIDE partner repository with the dataset identifier PXD031991. Code for reproducing the
1144 analysis is available in our lab Github repository ([https://github.com/sebepedroslab/chromatin-](https://github.com/sebepedroslab/chromatin-evolution-analysis)
1145 [evolution-analysis](https://github.com/sebepedroslab/chromatin-evolution-analysis)).

1146

1147

1148 **References**

- 1149 1. Struhl, K. Fundamentally different logic of gene regulation in eukaryotes and prokaryotes. *Cell* **98**, 1–
1150 4 (1999).
- 1151 2. Kornberg, R. D. & Lorch, Y. Primary Role of the Nucleosome. *Mol. Cell* **79**, 371–375 (2020).
- 1152 3. Jenuwein, T. & Allis, C. D. Translating the Histone Code. *Science* **293**, 1074–1080 (2001).
- 1153 4. Berger, S. L. The complex language of chromatin regulation during transcription. *Nature* **447**, 407–12
1154 (2007).
- 1155 5. Banaszynski, L. a, Allis, C. D. & Lewis, P. W. Histone variants in metazoan development. *Dev. Cell*
1156 **19**, 662–74 (2010).
- 1157 6. Allis, C. D. & Jenuwein, T. The molecular hallmarks of epigenetic control. *Nat. Rev. Genet.* **1**, (2016).
- 1158 7. Sultana, T. *et al.* The Landscape of L1 Retrotransposons in the Human Genome Is Shaped by Pre-
1159 insertion Sequence Biases and Post-insertion Selection. *Mol. Cell* **74**, 555-570.e7 (2019).
- 1160 8. Gangadharan, S., Mularoni, L., Fain-Thornton, J., Wheelan, S. J. & Craig, N. L. DNA transposon
1161 Hermes inserts into DNA in nucleosome-free regions in vivo. *Proc. Natl. Acad. Sci.* **107**, 21966–
1162 21972 (2010).
- 1163 9. Shinn, P. *et al.* HIV-1 Integration in the Human Genome Favors Active Genes and Local Hotspots.
1164 *Cell* **110**, 521–529 (2002).
- 1165 10. Goodier, J. L. Restricting retrotransposons: a review. *Mob. DNA* **7**, 16 (2016).

- 1166 11. Molaro, A. & Malik, H. S. Hide and seek: how chromatin-based pathways silence retroelements in the
1167 mammalian germline. *Curr. Opin. Genet. Dev.* **37**, 51–58 (2016).
- 1168 12. Malik, H. S. & Henikoff, S. Phylogenomics of the nucleosome. *Nat. Struct. Biol.* **10**, 882–91 (2003).
- 1169 13. Talbert, P. B. & Henikoff, S. Histone variants--ancient wrap artists of the epigenome. *Nat. Rev. Mol.*
1170 *Cell Biol.* **11**, 264–75 (2010).
- 1171 14. Soboleva, T. a., Nekrasov, M., Ryan, D. P. & Tremethick, D. J. Histone variants at the transcription
1172 start-site. *Trends Genet.* **30**, 199–209 (2014).
- 1173 15. Zink, L.-M. & Hake, S. B. Histone variants: nuclear function and disease. *Curr. Opin. Genet. Dev.* **37**,
1174 82–89 (2016).
- 1175 16. Weber, C. M. & Henikoff, S. Histone variants: dynamic punctuation in transcription. *Genes Dev.* **28**,
1176 672–82 (2014).
- 1177 17. Borg, M., Jiang, D. & Berger, F. Histone variants take center stage in shaping the epigenome. *Curr.*
1178 *Opin. Plant Biol.* **61**, 101991 (2021).
- 1179 18. Zentner, G. E. & Henikoff, S. Regulation of nucleosome dynamics by histone modifications. *Nat.*
1180 *Struct. Mol. Biol.* **20**, 259–66 (2013).
- 1181 19. Campos, E. I. & Reinberg, D. Histones: annotating chromatin. *Annu. Rev. Genet.* **43**, 559–99 (2009).
- 1182 20. Strahl, B. D. & Allis, C. D. The language of covalent histone modifications. *Nature* **403**, 41–45
1183 (2000).
- 1184 21. Bannister, A. J. & Kouzarides, T. Regulation of chromatin by histone modifications. *Cell Res.* **21**,
1185 381–95 (2011).
- 1186 22. Talbert, P. B. & Henikoff, S. The Yin and Yang of Histone Marks in Transcription. *Annu. Rev.*
1187 *Genomics Hum. Genet.* **22**, 147–170 (2021).
- 1188 23. Taverna, S. D., Li, H., Ruthenburg, A. J., Allis, C. D. & Patel, D. J. How chromatin-binding modules
1189 interpret histone modifications: lessons from professional pocket pickers. *Nat. Struct. Mol. Biol.* **14**,
1190 1025–40 (2007).
- 1191 24. Musselman, C. a, Lalonde, M.-E., Côté, J. & Kutateladze, T. G. Perceiving the epigenetic landscape
1192 through histone readers. *Nat. Struct. Mol. Biol.* **19**, 1218–27 (2012).
- 1193 25. Gurard-Levin, Z. a, Quivy, J.-P. & Almouzni, G. Histone Chaperones: Assisting Histone Traffic and
1194 Nucleosome Dynamics. *Annu. Rev. Biochem.* **83**, 487–517 (2014).
- 1195 26. Burgess, R. J. & Zhang, Z. Histone chaperones in nucleosome assembly and human disease. *Nat.*
1196 *Struct. Mol. Biol.* **20**, 14–22 (2013).
- 1197 27. Koster, M. J. E., Snel, B. & Timmers, H. T. M. Genesis of Chromatin and Transcription Dynamics in
1198 the Origin of Species. *Cell* **161**, 724–736 (2015).
- 1199 28. Hargreaves, D. C. & Crabtree, G. R. ATP-dependent chromatin remodeling: genetics, genomics and
1200 mechanisms. *Cell Res.* **21**, 396–420 (2011).
- 1201 29. Gornik, S. G. *et al.* Loss of nucleosomal DNA condensation coincides with appearance of a novel
1202 nuclear protein in dinoflagellates. *Curr. Biol.* **22**, 2303–12 (2012).
- 1203 30. Mattioli, F. *et al.* Structure of histone-based chromatin in Archaea. *Science* **357**, 609–612 (2017).
- 1204 31. Warnecke, T., Becker, E. a, Facciotti, M. T., Nislow, C. & Lehner, B. Conserved substitution patterns
1205 around nucleosome footprints in eukaryotes and Archaea derive from frequent nucleosome
1206 repositioning through evolution. *PLoS Comput. Biol.* **9**, e1003373 (2013).
- 1207 32. Ammar, R. *et al.* Chromatin is an ancient innovation conserved between Archaea and Eukarya. *Elife*
1208 **1**, e00078 (2012).

- 1209 33. Rojec, M., Hocher, A., Merckenschlager, M. & Warnecke, T. Chromatinization of Escherichia coli with
1210 archaeal histones. *bioRxiv* 660035 (2019) doi:10.1101/660035.
- 1211 34. Forbes, A. J. *et al.* Targeted analysis and discovery of posttranslational modifications in proteins from
1212 methanogenic archaea by top-down MS. *Proc. Natl. Acad. Sci. U. S. A.* **101**, 2678–83 (2004).
- 1213 35. Weidenbach, K. *et al.* Deletion of the archaeal histone in Methanosarcina mazei Gö1 results in
1214 reduced growth and genomic transcription. *Mol. Microbiol.* **67**, 662–671 (2008).
- 1215 36. Talbert, P. B., Meers, M. P. & Henikoff, S. Old cogs, new tricks: the evolution of gene expression in a
1216 chromatin context. *Nat. Rev. Genet.* (2019) doi:10.1038/s41576-019-0105-7.
- 1217 37. de Mendoza, A. & Sebe-Pedros, A. Origin and evolution of eukaryotic transcription factors. *Curr.*
1218 *Opin. Genet. Dev.* **59**, 25–32 (2019).
- 1219 38. Schwaiger, M. *et al.* Evolutionary conservation of the eumetazoan gene regulatory landscape.
1220 *Genome Res.* **24**, 639–650 (2014).
- 1221 39. Sebé-Pedrós, A. *et al.* Early metazoan cell type diversity and the evolution of multicellular gene
1222 regulation. *Nat. Ecol. Evol.* **2**, 1176–1188 (2018).
- 1223 40. Connolly, L. R., Smith, K. M. & Freitag, M. The Fusarium graminearum Histone H3 K27
1224 Methyltransferase KMT6 Regulates Development and Expression of Secondary Metabolite Gene
1225 Clusters. *PLoS Genet.* **9**, e1003916 (2013).
- 1226 41. Jamieson, K., Rountree, M. R., Lewis, Z. a, Stajich, J. E. & Selker, E. U. Regional control of histone
1227 H3 lysine 27 methylation in Neurospora. *Proc. Natl. Acad. Sci.* **110**, 6027–6032 (2013).
- 1228 42. Sebé-Pedrós, A. *et al.* The Dynamic Regulatory Genome of Capsaspora and the Origin of Animal
1229 Multicellularity. *Cell* **165**, 1224–1237 (2016).
- 1230 43. Bourdareau, S. *et al.* Histone modifications during the life cycle of the brown alga Ectocarpus.
1231 *Genome Biol.* **22**, 12 (2021).
- 1232 44. Wang, S. Y. *et al.* Role of epigenetics in unicellular to multicellular transition in Dictyostelium.
1233 *Genome Biol.* **22**, 134 (2021).
- 1234 45. Taverna, S. D., Coyne, R. S. & Allis, C. D. Methylation of Histone H3 at Lysine 9 Targets
1235 Programmed DNA Elimination in Tetrahymena University of Virginia Health System. *Cell* **110**, 701–
1236 711 (2002).
- 1237 46. Garcia, B. a *et al.* Organismal differences in post-translational modifications in histones H3 and H4.
1238 *J. Biol. Chem.* **282**, 7641–55 (2007).
- 1239 47. Drinnenberg, I. A. *et al.* EvoChromo: towards a synthesis of chromatin biology and evolution.
1240 *Development* **146**, dev178962 (2019).
- 1241 48. Draizen, E. J. *et al.* HistoneDB 2.0: a histone database with variants—an integrated resource to
1242 explore histones and their variants. *Database* **2016**, baw014 (2016).
- 1243 49. Maile, T. M. *et al.* Mass Spectrometric Quantification of Histone Post-translational Modifications by
1244 a Hybrid Chemical Labeling Method. *Mol. Cell. Proteomics* **14**, 1148–1158 (2015).
- 1245 50. Li, B., Carey, M. & Workman, J. L. The Role of Chromatin during Transcription. *Cell* **128**, 707–719
1246 (2007).
- 1247 51. Rajagopal, N. *et al.* Distinct and Predictive Histone Lysine Acetylation Patterns at Promoters,
1248 Enhancers, and Gene Bodies. *G3 Genes/Genomes/Genetics* **4**, 2051–2063 (2014).
- 1249 52. Koonin, E. V. & Yutin, N. The Dispersed Archaeal Eukaryome and the Complex Archaeal Ancestor of
1250 Eukaryotes. *Cold Spring Harb. Perspect. Biol.* **6**, a016188–a016188 (2014).
- 1251 53. Sandman, K. & Reeve, J. N. Archaeal histones and the origin of the histone fold. *Curr. Opin.*
1252 *Microbiol.* **9**, 520–5 (2006).

- 1253 54. Pereira, S. L., Grayling, R. a, Lurz, R. & Reeve, J. N. Archaeal nucleosomes. *Proc. Natl. Acad. Sci.*
1254 *U. S. A.* **94**, 12633–7 (1997).
- 1255 55. Henneman, B., van Emmerik, C., van Ingen, H. & Dame, R. T. Structure and function of archaeal
1256 histones. *PLOS Genet.* **14**, e1007582 (2018).
- 1257 56. Imachi, H. *et al.* Isolation of an archaeon at the prokaryote-eukaryote interface. *bioRxiv* 726976
1258 (2019) doi:10.1101/726976.
- 1259 57. Spang, A. *et al.* Complex archaea that bridge the gap between prokaryotes and eukaryotes. *Nature*
1260 **521**, 173–179 (2015).
- 1261 58. Zaremba-Niedzwiedzka, K. *et al.* Asgard archaea illuminate the origin of eukaryotic cellular
1262 complexity. *Nature* **541**, 353–358 (2017).
- 1263 59. Da Cunha, V., Gaia, M., Nasir, A. & Forterre, P. Asgard archaea do not close the debate about the
1264 universal tree of life topology. *PLOS Genet.* **14**, e1007215 (2018).
- 1265 60. Alva, V. & Lupas, A. N. Histones predate the split between bacteria and archaea. *Bioinformatics* **35**,
1266 2349–2353 (2019).
- 1267 61. Allis, C. D. *et al.* New Nomenclature for Chromatin-Modifying Enzymes. *Cell* **131**, 633–636 (2007).
- 1268 62. Wu, F. *et al.* Unique mobile elements and scalable gene flow at the prokaryote–eukaryote boundary
1269 revealed by circularized Asgard archaea genomes. *Nat. Microbiol.* **7**, 200–212 (2022).
- 1270 63. Schuettengruber, B., Bourbon, H.-M., Di Croce, L. & Cavalli, G. Genome Regulation by Polycomb
1271 and Trithorax: 70 Years and Counting. *Cell* **171**, 34–57 (2017).
- 1272 64. Dion, M. F., Altschuler, S. J., Wu, L. F. & Rando, O. J. Genomic characterization reveals a simple
1273 histone H4 acetylation code. *Proc. Natl. Acad. Sci.* **102**, 5501–5506 (2005).
- 1274 65. de Jong, J. *et al.* Chromatin Landscapes of Retroviral and Transposon Integration Profiles. *PLoS*
1275 *Genet.* **10**, e1004250 (2014).
- 1276 66. Sultana, T., Zamborlini, A., Cristofari, G. & Lesage, P. Integration site selection by retroviruses and
1277 transposable elements in eukaryotes. *Nat. Rev. Genet.* **18**, 292–308 (2017).
- 1278 67. Gao, X., Hou, Y., Ebina, H., Levin, H. L. & Voytas, D. F. Chromodomains direct integration of
1279 retrotransposons to heterochromatin. *Genome Res.* **18**, 359–369 (2008).
- 1280 68. Cosby, R. L. *et al.* Recurrent evolution of vertebrate transcription factors by transposase capture.
1281 *Science* **371**, eabc6405 (2021).
- 1282 69. Cordaux, R., Udit, S., Batzer, M. A. & Feschotte, C. Birth of a chimeric primate gene by capture of
1283 the transposase gene from a mobile element. *Proc. Natl. Acad. Sci.* **103**, 8101 LP – 8106 (2006).
- 1284 70. Fiedler, M. *et al.* Decoding of Methylated Histone H3 Tail by the Pygo-BCL9 Wnt Signaling
1285 Complex. *Mol. Cell* **30**, 507–518 (2008).
- 1286 71. Erives, A. J. Phylogenetic analysis of the core histone doublet and DNA topo II genes of
1287 *Marseilleviridae*: evidence of proto-eukaryotic provenance. *Epigenetics Chromatin* **10**, 55 (2017).
- 1288 72. Liu, Y. *et al.* Virus-encoded histone doublets are essential and form nucleosome-like structures. *Cell*
1289 **184**, 4237–4250.e19 (2021).
- 1290 73. Valencia-Sánchez, M. I. *et al.* The structure of a virus-encoded nucleosome. *Nat. Struct. Mol. Biol.*
1291 **28**, 413–417 (2021).
- 1292 74. Iyer, L. M., Balaji, S., Koonin, E. V & Aravind, L. Evolutionary genomics of nucleo-cytoplasmic
1293 large DNA viruses. *Virus Res.* **117**, 156–184 (2006).
- 1294 75. Nagamine, T. Apoptotic arms races in insect–baculovirus coevolution. *Physiol. Entomol.* phen.12371
1295 (2021) doi:10.1111/phen.12371.

- 1296 76. Starrett, G. J. *et al.* Adintoviruses: An Animal-Tropic Family of Midsize Eukaryotic Linear dsDNA
1297 (MELD) Viruses. *bioRxiv* 697771 (2020) doi:10.1101/697771.
- 1298 77. Hocher, A. *et al.* Growth temperature is the principal driver of chromatinization in archaea. *bioRxiv*
1299 2021.07.08.451601 (2021) doi:10.1101/2021.07.08.451601.
- 1300 78. Alpha-Bazin, B. *et al.* Lysine-specific acetylated proteome from the archaeon *Thermococcus*
1301 *gammatolerans* reveals the presence of acetylated histones. *J. Proteomics* **232**, 104044 (2021).
- 1302 79. Eme, L., Spang, A., Lombard, J., Stairs, C. W. & Ettema, T. J. G. Archaea and the origin of
1303 eukaryotes. *Nat. Rev. Microbiol.* **15**, 711–723 (2017).
- 1304 80. Ak1l, C. & Robinson, R. C. Genomes of Asgard archaea encode profilins that regulate actin. *Nature*
1305 **562**, 439–443 (2018).
- 1306 81. Koonin, E. V. The origin and early evolution of eukaryotes in the light of phylogenomics. *Genome*
1307 *Biol.* **11**, 209 (2010).
- 1308 82. Sebé-Pedrós, A., Grau-Bové, X., Richards, T. a & Ruiz-Trillo, I. Evolution and Classification of
1309 Myosins, a Paneukaryotic Whole-Genome Approach. *Genome Biol. Evol.* **6**, 290–305 (2014).
- 1310 83. Richards, T. A. & Cavalier-Smith, T. Myosin domain evolution and the primary divergence of
1311 eukaryotes. *Nature* **436**, 1113–8 (2005).
- 1312 84. Wickstead, B., Gull, K. & Richards, T. Patterns of kinesin evolution reveal a complex ancestral
1313 eukaryote with a multifunctional cytoskeleton. *BMC Evol. Biol.* **10**, 110 (2010).
- 1314 85. Dacks, J. B. & Field, M. C. Evolution of the eukaryotic membrane-trafficking system: origin, tempo
1315 and mode. *J. Cell Sci.* **120**, 2977–85 (2007).
- 1316 86. Collins, L. & Penny, D. Complex Spliceosomal Organization Ancestral to Extant Eukaryotes. *Mol.*
1317 *Biol. Evol.* **22**, 1053–1066 (2005).
- 1318 87. Grau-Bové, X., Sebé-Pedrós, A. & Ruiz-Trillo, I. The Eukaryotic Ancestor Had a Complex Ubiquitin
1319 Signaling System of Archaeal Origin. *Mol. Biol. Evol.* **32**, 726–739 (2015).
- 1320 88. Kundaje, A. *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317–330
1321 (2015).
- 1322 89. Ho, J. W. K. *et al.* Comparative analysis of metazoan chromatin organization. *Nature* **512**, 449–452
1323 (2014).
- 1324 90. Montgomery, S. A. *et al.* Chromatin Organization in Early Land Plants Reveals an Ancestral
1325 Association between H3K27me3, Transposons, and Constitutive Heterochromatin. *Curr. Biol.* **30**,
1326 573-588.e7 (2020).
- 1327 91. Frapporti, A. *et al.* The Polycomb protein Ez11 mediates H3K9 and H3K27 methylation to repress
1328 transposable elements in *Paramecium*. *Nat. Commun.* **10**, 2710 (2019).
- 1329 92. Lennartsson, A. & Ekwall, K. Histone modification patterns and epigenetic codes. *Biochim. Biophys.*
1330 *Acta* **1790**, 863–8 (2009).
- 1331 93. Peterson, C. L. & Laniel, M.-A. Histones and histone modifications. *Curr. Biol.* **14**, R546-51 (2004).
- 1332 94. Rando, O. J. Combinatorial complexity in chromatin structure and function: revisiting the histone
1333 code. *Curr. Opin. Genet. Dev.* **22**, 148–155 (2012).
- 1334 95. de Mendoza, A., Pflueger, J. & Lister, R. Capture of a functionally active methyl-CpG binding
1335 domain by an arthropod retrotransposon family. *Genome Res.* **29**, 1277–1286 (2019).
- 1336 96. De Mendoza, A. *et al.* Recurrent acquisition of cytosine methyltransferases into eukaryotic
1337 retrotransposons. *Nat. Commun.* **9**, 1–11 (2018).
- 1338 97. Ji, X. *et al.* Chromatin proteomic profiling reveals novel proteins associated with histone-marked
1339 genomic regions. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 3841–3846 (2015).

- 1340 98. Wierer, M. & Mann, M. Proteomics to study DNA-bound and chromatin-associated gene regulatory
1341 complexes. *Hum. Mol. Genet.* **25**, R106–R114 (2016).
- 1342 99. Villaseñor, R. *et al.* ChromID identifies the protein interactome at chromatin marks. *Nat. Biotechnol.*
1343 **38**, 728–736 (2020).
- 1344 100. Stieglmeier, M. *et al.* Nitrososphaera viennensis gen. nov., sp. nov., an aerobic and mesophilic,
1345 ammonia-oxidizing archaeon from soil and a member of the archaeal phylum Thaumarchaeota. *Int. J.*
1346 *Syst. Evol. Microbiol.* **64**, 2738–2752 (2014).
- 1347 101. Tirichine, L. *et al.* Histone extraction protocol from the two model diatoms Phaeodactylum
1348 tricornutum and Thalassiosira pseudonana. *Mar. Genomics* **13**, 21–25 (2014).
- 1349 102. Shechter, D., Dormann, H. L., Allis, C. D. & Hake, S. B. Extraction, purification and analysis of
1350 histones. *Nat. Protoc.* **2**, 1445–57 (2007).
- 1351 103. Perkins, D. N., Pappin, D. J. C., Creasy, D. M. & Cottrell, J. S. Probability-based protein
1352 identification by searching sequence databases using mass spectrometry data. *Electrophoresis* **20**,
1353 3551–3567 (1999).
- 1354 104. Taus, T. *et al.* Universal and confident phosphorylation site localization using phosphoRS. *J.*
1355 *Proteome Res.* **10**, 5354–5362 (2011).
- 1356 105. Vizcaíno, J. A. *et al.* 2016 update of the PRIDE database and its related tools. *Nucleic Acids Res.* **44**,
1357 D447-56 (2016).
- 1358 106. Hagberg, A. A., Schult, D. A. & Swart, P. J. Exploring Network Structure, Dynamics, and Function
1359 using NetworkX. in *Proceedings of the 7th Python in Science Conference* (eds. Varoquaux, G.,
1360 Vaught, T. & Millman, J.) 11–15 (2008).
- 1361 107. Katoh, K. & Standley, D. M. MAFFT Multiple Sequence Alignment Software Version 7:
1362 Improvements in Performance and Usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).
- 1363 108. Nguyen, L.-T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: A Fast and Effective
1364 Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol. Biol. Evol.* **32**, 268–274
1365 (2015).
- 1366 109. Veluchamy, A. *et al.* An integrative analysis of post-translational histone modifications in the marine
1367 diatom Phaeodactylum tricornutum. *Genome Biol.* **16**, 102 (2015).
- 1368 110. Ren, Q. & Gorovsky, M. a. Histone H2A.Z acetylation modulates an essential charge patch. *Mol. Cell*
1369 **7**, 1329–35 (2001).
- 1370 111. Allis, C. D. *et al.* hv1 is an evolutionarily conserved H2A variant that is preferentially associated with
1371 active genes. *J. Biol. Chem.* **261**, 1941–1948 (1986).
- 1372 112. Fusauchi, Y. & Iwai, K. Tetrahymena histone H2A. Acetylation in the N-terminal sequence and
1373 phosphorylation in the C-terminal sequence. *J. Biochem.* **95**, 147–154 (1984).
- 1374 113. Xiong, L., Adhvaryu, K. K., Selker, E. U. & Wang, Y. Mapping of Lysine Methylation and
1375 Acetylation in Core Histones of Neurospora crassa. *Biochemistry* **49**, 5236–5243 (2010).
- 1376 114. Zhang, K., Sridhar, V. V., Zhu, J., Kapoor, A. & Zhu, J. K. Distinctive core histone post-translational
1377 modification patterns in Arabidopsis thaliana. *PLoS One* **2**, (2007).
- 1378 115. Johnson, L. *et al.* Mass spectrometry analysis of Arabidopsis histone H3 reveals distinct combinations
1379 of post-translational modifications. *Nucleic Acids Res.* **32**, 6511–6518 (2004).
- 1380 116. Bergmüller, E., Gehrig, P. M. & Gruissem, W. Characterization of post-translational modifications of
1381 histone H2B-variants isolated from Arabidopsis thaliana. *J. Proteome Res.* **6**, 3655–3668 (2007).
- 1382 117. Beck, H. C. *et al.* Quantitative Proteomic Analysis of Post-translational Modifications of Human
1383 Histones. *Mol. Cell. Proteomics* **5**, 1314–1325 (2006).

- 1384 118. Goudarzi, A. *et al.* Dynamic Competing Histone H4 K5K8 Acetylation and Butyrylation Are
1385 Hallmarks of Highly Active Gene Promoters. *Mol. Cell* **62**, 169–180 (2016).
- 1386 119. Hake, S. B. *et al.* Expression patterns and post-translational modifications associated with
1387 mammalian histone H3 variants. *J. Biol. Chem.* **281**, 559–568 (2006).
- 1388 120. Tan, M. *et al.* Identification of 67 histone marks and histone lysine crotonylation as a new type of
1389 histone modification. *Cell* **146**, 1016–1028 (2011).
- 1390 121. Moniruzzaman, M., Martinez-Gutierrez, C. A., Weinheimer, A. R. & Aylward, F. O. Dynamic genome
1391 evolution and complex virocell metabolism of globally-distributed giant viruses. *Nat. Commun.* **11**,
1392 1–11 (2020).
- 1393 122. Punta, M. *et al.* The Pfam protein families database. *Nucleic Acids Res.* **40**, D290-301 (2012).
- 1394 123. Eddy, S. R. Accelerated profile HMM searches. *PLoS Comput. Biol.* **7**, e1002195 (2011).
- 1395 124. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND. *Nat.*
1396 *Methods* **12**, 59–60 (2015).
- 1397 125. Enright, A. J., Van Dongen, S. & Ouzounis, C. A. An efficient algorithm for large-scale detection of
1398 protein families. *Nucleic Acids Res.* **30**, 1575–1584 (2002).
- 1399 126. Steenwyk, J. L., Buida, T. J., Li, Y., Shen, X.-X. & Rokas, A. ClipKIT: A multiple sequence alignment
1400 trimming software for accurate phylogenomic inference. *PLoS Biol.* **18**, e3001007 (2020).
- 1401 127. Minh, B. Q., Nguyen, M. A. T. & von Haeseler, A. Ultrafast approximation for phylogenetic
1402 bootstrap. *Mol. Biol. Evol.* **30**, 1188–95 (2013).
- 1403 128. Grau-Bové, X. & Sebé-Pedrós, A. Orthology Clusters from Gene Trees with Possvm. *Mol. Biol. Evol.*
1404 **38**, 5204–5208 (2021).
- 1405 129. Huerta-Cepas, J., Dopazo, H., Dopazo, J. & Gabaldón, T. The human phylome. *Genome Biol.* **8**, R109
1406 (2007).
- 1407 130. Huerta-Cepas, J., Serra, F. & Bork, P. ETE 3: Reconstruction, Analysis, and Visualization of
1408 Phylogenomic Data. *Mol. Biol. Evol.* **33**, 1635–1638 (2016).
- 1409 131. Csűrös, M. & Miklós, I. A Probabilistic Model for Gene Content Evolution with Duplication, Loss,
1410 and Horizontal Transfer BT - Research in Computational Molecular Biology. in (eds. Apostolico, A.,
1411 Guerra, C., Istrail, S., Pevzner, P. A. & Waterman, M.) 206–220 (Springer Berlin Heidelberg, 2006).
- 1412 132. Csűrös, M. Count: evolutionary analysis of phylogenetic profiles with parsimony and likelihood.
1413 *Bioinformatics* **26**, 1910–2 (2010).
- 1414 133. Gansner, E. R. & North, S. C. An open graph visualization system and its applications to software
1415 engineering. *Softw. Pract. Exp.* **30**, 1203–1233 (2000).
- 1416 134. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-generation
1417 sequencing data. *Bioinformatics* **28**, 3150–3152 (2012).
- 1418 135. Jombart, T., Balloux, F. & Dray, S. adephylo: new tools for investigating the phylogenetic signal in
1419 biological traits. *Bioinformatics* **26**, 1907–1909 (2010).
- 1420 136. Wells, J. N. & Feschotte, C. A Field Guide to Eukaryotic Transposable Elements. *Annu. Rev. Genet.*
1421 **54**, annurev-genet-040620-022145 (2020).
- 1422 137. Storer, J., Hubley, R., Rosen, J., Wheeler, T. J. & Smit, A. F. The Dfam community resource of
1423 transposable element families, sequence models, and genome annotations. *Mob. DNA* **12**, 2 (2021).
- 1424 138. Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421 (2009).
- 1425 139. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform.
1426 *Bioinformatics* **25**, 1754–1760 (2009).

- 1427 140. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features.
1428 *Bioinforma.* **26**, 841–842 (2010).
- 1429 141. Anisimova, M., Gil, M., Dufayard, J.-F., Dessimoz, C. & Gascuel, O. Survey of Branch Support
1430 Methods Demonstrates Accuracy, Power, and Robustness of Fast Likelihood-based Approximation
1431 Schemes. *Syst. Biol.* **60**, 685–699 (2011).
- 1432 142. Bodenhofer, U., Bonatesta, E., Horejš-Kainrath, C. & Hochreiter, S. msa: an R package for multiple
1433 sequence alignment. *Bioinformatics* **31**, 3997–3999 (2015).
- 1434
- 1435

Figure 1

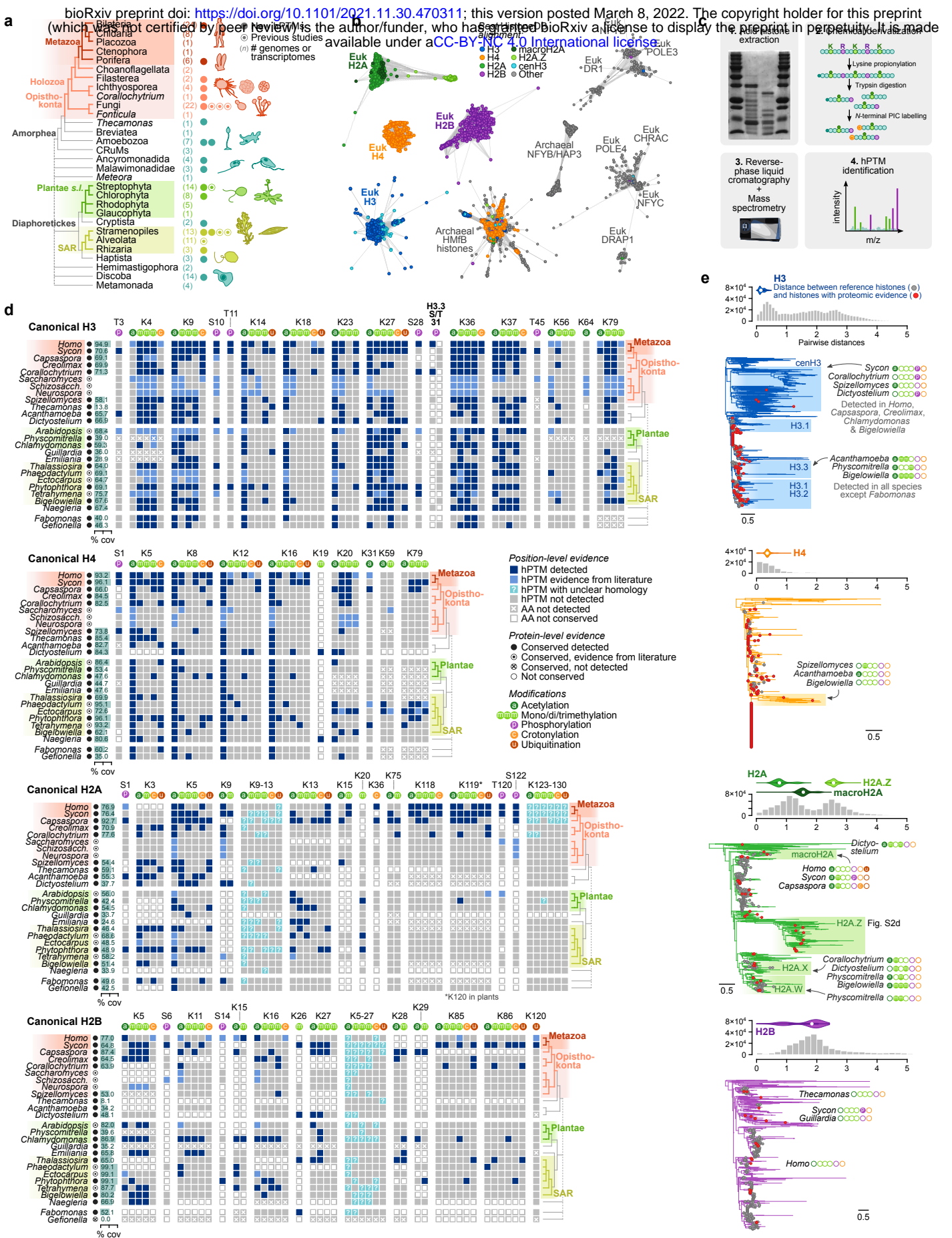


Figure 2

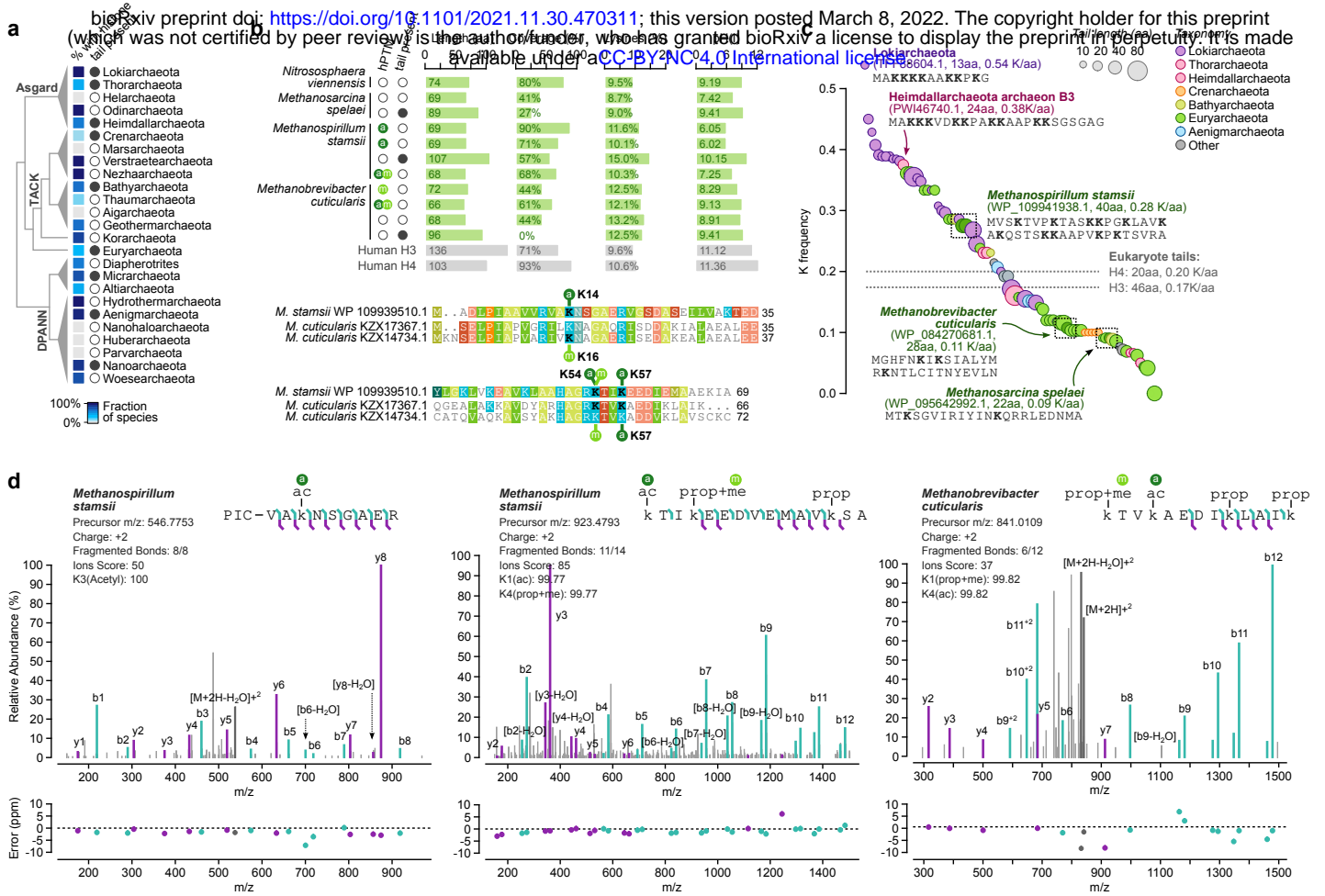


Figure 3

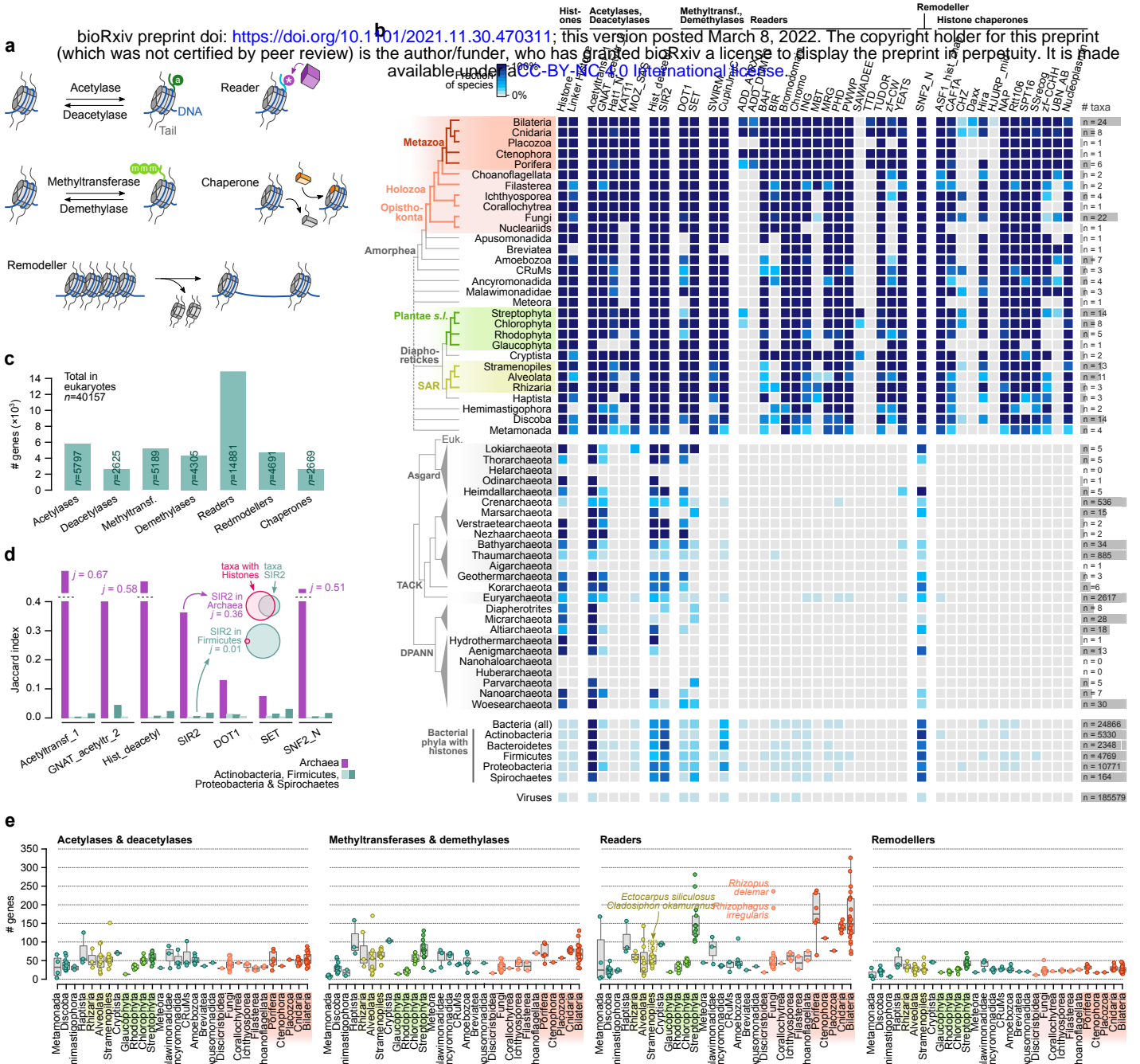


Figure 4

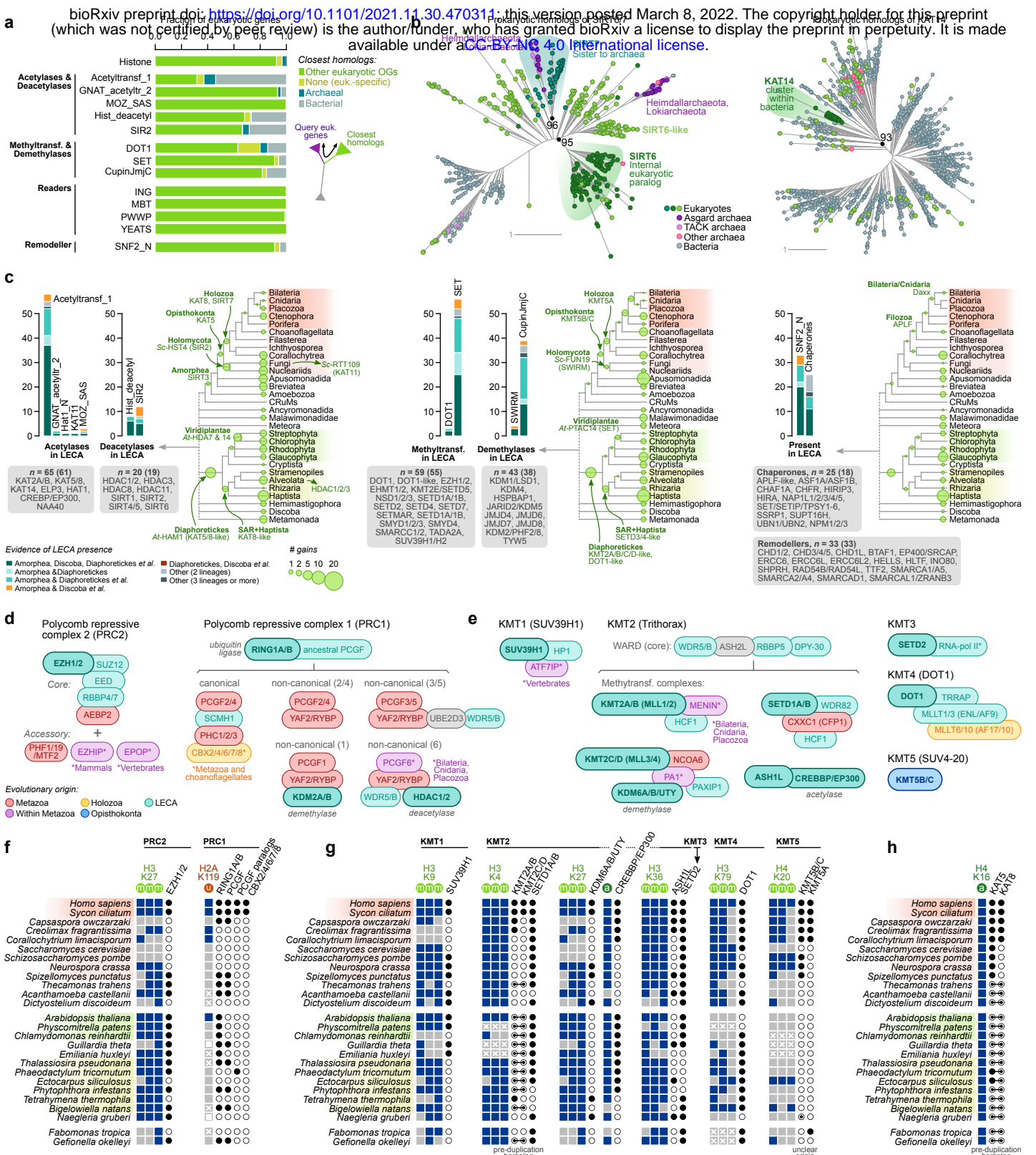


Figure 5

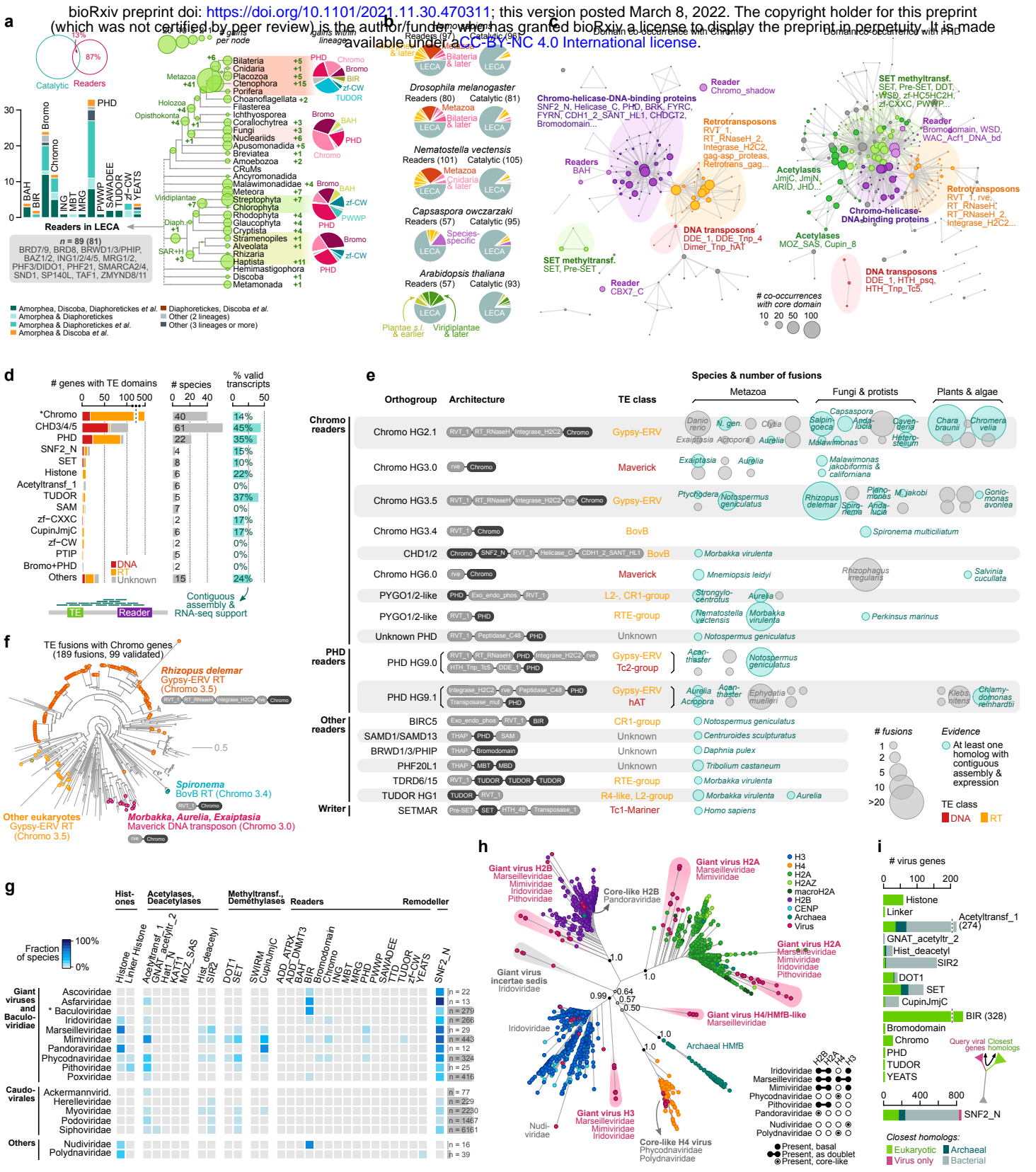


Figure 6

bioRxiv preprint doi: <https://doi.org/10.1101/2021.11.30.470311>; this version posted March 8, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC 4.0 International license.

