

1
2
3 **1 Deep conservation of histone variants in Thermococcales archaea**
4
5 **2**

6 **3 Kathryn M Stevens^{1,2}, Antoine Hocher^{1,2}, Tobias Warnecke^{1,2*}**
7
8 **4**

9 **5 ¹Medical Research Council London Institute of Medical Sciences, London, United Kingdom**

10 **6 ²Institute of Clinical Sciences, Faculty of Medicine, Imperial College London, London,**
11 **7 United Kingdom**

12 **8**
13 **9 *corresponding author: tobias.warnecke@lms.mrc.ac.uk**
14
15
16
17
18
19 **10**
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1 Abstract

2
3
4
5
6
7 Histones are ubiquitous in eukaryotes where they assemble into nucleosomes, binding and
8 wrapping DNA to form chromatin. One process to modify chromatin and regulate DNA
9 accessibility is the replacement of histones in the nucleosome with paralogous variants.
10 Histones are also present in archaea but whether and how histone variants contribute to the
11 generation of different physiologically relevant chromatin states in these organisms remains
12 largely unknown. Conservation of paralogs with distinct properties can provide *prima facie*
13 evidence for defined functional roles. We recently revealed deep conservation of histone
14 paralogs with different properties in the Methanobacteriales, but little is known
15 experimentally about these histones. In contrast, the two histones of the model archaeon
16 *Thermococcus kodakarensis*, HTkA and HTkB, have been examined in some depth, both *in*
17 *vitro* and *in vivo*. HTkA and HTkB exhibit distinct DNA-binding behaviours and elicit unique
18 transcriptional responses when deleted. Here, we consider the evolution of HTkA/B and their
19 orthologs across the order Thermococcales. We find histones with signature HTkA- and
20 HTkB-like properties to be present in almost all Thermococcales genomes. Phylogenetic
21 analysis indicates the presence of one HTkA- and one HTkB-like histone in the ancestor of
22 Thermococcales and long-term maintenance of these two paralogs throughout
23 Thermococcales diversification. Our results support the notion that archaea and eukaryotes
24 have convergently evolved histone variants that carry out distinct adaptive functions.
25 Intriguingly, we also detect more highly diverged histone-fold proteins, related to those found
26 in some bacteria, in several Thermococcales genomes. The functions of these bacteria-type
27 histones remain unknown, but structural modelling suggests that they can form heterodimers
28 with HTkA/B-like histones.
29
30
31
32
33
34

1
2
3 **1 Significance statement**
4
5 **2**

6
7 **3** Histones are key components of chromatin in eukaryotes and many archaea. In eukaryotes,
8 **4** histone variants exist that play defined roles in cellular function and development. Some of
9 **5** these variants are highly conserved and can date back to the last common ancestor of
10 **6** eukaryotes. Archaea also frequently encode multiple sequence-divergent histones but whether
11 **7** these play distinct functional roles that are conserved through evolution remains largely
12 **8** unknown. Here, using phylogenetic tools, we establish the existence of histone variants with
13 **9** different properties in the order Thermococcales, which have been conserved for hundreds of
14 **10** millions of years. Our work provides additional evidence for histone variants in archaea and
15 **11** that these have evolved in parallel to eukaryotes to mediate flexible, adaptive chromatin
16 **12** states.
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1 Introduction

2
3 The ability of eukaryotic cells to respond to environmental change and regulate transcription
4 relies on dynamic control of DNA accessibility through chromatin alterations. This involves
5 many different processes, including the addition/removal of histone modifications and the
6 exchange of histone proteins for paralogous variants. Such variants can modify structural
7 properties of the nucleosome or change how it interacts with its binding partners (Talbert &
8 Henikoff 2010; Martire & Banaszynski 2020; Henikoff & Smith 2015). For example,
9 macroH2A has a large C-terminal domain and precipitates transcriptional repression (Martire
10 & Banaszynski 2020; Bönisch & Hake 2012) while cenH3, a fast-evolving H3 variant, is
11 specifically localised to centromeres and involved in chromosome segregation (Palmer et al.
12 1991; Talbert & Henikoff 2010). Importantly, significant functional changes can come from
13 small differences in sequence. H3.3, for example, is deposited in a replication-independent
14 manner in actively transcribed regions of the genome (Talbert & Henikoff 2010) and
15 important for mammalian development (Sitbon et al. 2020; Jang et al. 2015), but differs from
16 its paralog H3.1 by only five amino acid.

17
18 Histones are not exclusive to eukaryotes. Archaeal histone proteins, first discovered in
19 *Methanothermus fervidus* (Starich et al. 1996; Sandman et al. 1990), have since been
20 identified in diverse archaeal lineages (Henneman et al. 2018; Hocher et al. 2021) and are
21 often highly expressed (Hocher et al. 2021). Eukaryotic and archaeal histones share a
22 conserved histone fold (HF) domain, form dimers and tetramers that are structurally very
23 similar, and bind DNA non-specifically, albeit with a preference for more bendable
24 sequences (Rojec et al. 2019; Nalabothula et al. 2013; Bailey et al. 2000; Decanniere et al.
25 2000; Mattioli et al. 2017). Unlike eukaryotic histones, almost all archaeal histones lack long
26 terminal extensions (“tails”) (Henneman et al. 2018). In at least some instances, archaeal
27 histones have the capacity to form homo- as well as heterodimers and to assemble into long
28 oligomeric structures on DNA (Mattioli et al. 2017). These extended complexes (Bowerman
29 et al. 2021) can, in theory, consist of different histone paralogs, providing opportunities for
30 chromatin state modulation through the exchange of histones with different properties
31 (Stevens et al. 2020). In fact, many archaea encode two or more sequence-divergent histone
32 paralogs, but whether these paralogs have defined functional roles akin to eukaryotic histone
33 variants, and whether their expression and assembly change dynamically to mediate adaptive
34 chromatin states, remains poorly understood.

1
2
3 1
4
5 2 What we do know from prior experimental work is that archaeal histone paralogs are more
6
7 3 than mere copy number variants. The two histones of *M. fervidus* (HMfA, HMfB), for
8
9 4 example, display differences in DNA binding affinity (Bailey et al. 2002). Compared to
10
11 5 HMfA, recombinant HMfB also induces more positive supercoiling upon binding to plasmid
12
13 6 DNA and forms a more compact complex as inferred from gel-shift and tethered particle
14
15 7 motion experiments (Sandman et al. 1994; Henneman et al. 2021). There are also differences
16
17 8 between HMfA and HMfB in their relative expression during the growth cycle: in early
18
19 9 exponential phase, HMfA is more highly expressed than HMfB, expression of which
20
21 10 increases towards stationary phase to reach an almost equal ratio between the two (Sandman
22
23 11 et al. 1994). The different properties of *M. fervidus* histones are consistent with the
24
25 12 hypothesis that the two paralogs may have distinct functions in nucleoid/chromatin biology,
26
27 13 but whether the properties are physiologically relevant and affect organismal fitness remains
28
29 14 to be addressed experimentally.
30
31 15

32
33 16 Recently, we considered this question using an evolutionary approach. We identified histone
34
35 17 paralogs in the order Methanobacteriales that exhibit distinct structural properties and have
36
37 18 been maintained over hundreds of millions of years (Stevens et al. 2020), indicative of the
38
39 19 importance of each individual paralog for fitness. Structural modelling identified histone
40
41 20 variants that prevent stable tetramerization and might act as *capstones* that limit further
42
43 21 extension when incorporated into a histone oligomer, providing a potential pathway to
44
45 22 dynamically alter chromatin state. Are the Methanobacteriales unique or are there other
46
47 23 clades of archaea with histone paralogs that have been maintained over long periods of time?
48
49 24 And do these paralogs also show conserved and distinct structural properties?
50
51 25

52
53 26 Here, we consider archaea in the order Thermococcales, which includes the model archaea
54
55 27 *Pyrococcus furiosus* and *Pyrococcus abyssi* as well as *T. kodakarensis*, which has served as a
56
57 28 model species for the *in vivo* study of archaeal histones. Thanks to the efforts of Santangelo
58
59 29 and co-workers in particular, its two histones – HTkA (TK1413) and HTkB (TK2289) – are
60
30 arguably the best characterized paralogs *in vivo*. Similar to HMfA and HMfB in *M. fervidus*,
31
32 31 HTkA and HTkB can assemble into long oligomers both *in vitro* and *in vivo* (Mattioli et al.
33
34 32 2017; Sanders et al. 2021; Bowerman et al. 2021; Maruyama et al. 2013). The two histones
35
36 33 differ from one another at 11 out of 67 residues (84% identity) and have several distinct
37
38 34 properties. HTkA is the more highly expressed paralog, at least in exponential phase, where it

1 makes up 1.1% of the proteome compared to 0.66% for HTkB (Hoher et al. 2021). Together,
2 they are abundant enough to coat the entire *T. kodakarensis* genome (Sanders et al. 2019).
3 HTkB has been shown to bind to DNA more strongly than HTkA and to form more compact
4 complexes, which show faster migration during agarose gel electrophoresis (Higashibata et
5 al. 1999). Deletion of each histone individually results in overlapping but distinct
6 perturbations of the transcriptome (Čuboňová et al. 2012; Sanders et al. 2021). Notably,
7 HTkB-deficient cells exhibit reduced growth, possibly due to changes in the expression, not
8 seen in strains lacking HTkA, of genes that encode translation factors and ribosomal proteins
9 (Čuboňová et al. 2012). Deletion of *htkA* but not *htkB* leads to downregulation of hypothetical
10 membrane proteins and prevents transformation of *T. kodakarensis*, suggesting HTkA alone
11 plays a critical role in DNA uptake and/or integration (Čuboňová et al. 2012).

12
13 In this study, we show that histone paralogs with HTkA- and HTkB-like properties are
14 present across the Thermococcales, including *Thermococcus*, *Pyrococcus*, and *Palaeococcus*
15 *spp.*. We use structural modelling to show that, in most Thermococcales, HTkB-like histones
16 are predicted to exhibit stronger DNA binding than those with HTkA-like properties.
17 Phylogenetic analysis reveals that HTkA-like histones share a common ancestor to the
18 exclusion of HTkB-like histones and *vice versa*, suggesting that the last common ancestor of
19 the Thermococcales already encoded an HTkA-like and an HTkB-like histone, each of which
20 has been maintained throughout the diversification of this clade for (very) approximately 750
21 million years (Wolfe & Fournier 2018). The long-term preservation of these two paralogs
22 across the order Thermococcales supports the notion that HTkA/B in *T. kodakarensis* (and
23 their orthologs in other Thermococcales) make unique contributions to genome function and
24 fitness. These findings add further evidence that histone variants are widespread in archaea,
25 evolving in parallel to those in eukaryotes.

26
27 Intriguingly, many Thermococcales archaea encode additional types of histone-fold proteins
28 that are similar to histone-fold proteins found in some bacteria (Alva & Lupas 2019). One of
29 these consists of an end-to-end duplication of the histone fold and is rarely found in archaea
30 outside the Thermococcales. Their physiological roles remain unknown, but structural
31 modelling suggest that they are able to heterodimerize with HTkA/B and might therefore
32 further diversify histone-based chromatin states in these archaea.

1 Results and Discussion

2
3
4
5
6
7 *Almost all Thermococcales have histone paralogs with HTkA- and HTkB-like properties*

8
9
10 To identify putative histone proteins across the Thermococcales, we scanned 61 predicted
11 proteomes using HMM models and iterative jackhmmer searches (see Methods). Histones
12 with a single histone-fold domain (similar to archaeal HMf-like histones) were found in all
13 genomes (Fig 1c, Table S1). We also recovered putative histone-fold proteins similar to those
14 found in some bacteria (Alva & Lupas 2019), which we discuss further below. A principal
15 component analysis of the HMf-like archaeal histones based on their amino acid properties
16 and isoelectric points (see Methods) suggests that histones can be assigned to one of two
17 groups. One of these groups contains HTkA, the other HTkB (Fig 1a). This is consistent with
18 a previous classification effort that also recovered two major groups of Thermococcales
19 histones (Henneman 2019). Amino acid identities at several residues along the histone fold
20 differ systematically between groups and are diagnostic of group membership. For example,
21 tyrosine is always found at position 35 (Y35) in HTkA-like histones whereas HTkB-like
22 histones have a positively charged lysine (K, 60 out of 62) or histidine (H, 2 out of 62).
23 Similarly, glutamic acid at position 18 (E18) is present in 59 out of the 61 HTkA-like but
24 none of the HTkB-like histones (Fig 1b).

25
26
27 For some of these residues, we know from prior *in vitro* studies – as well as structural
28 modelling – that amino acid identity can affect specific histone properties (Stevens et al.
29 2020; Soares et al. 2000). For example, substituting Y for K at residue 35 (the amino acids
30 seen in HTKA- and HTKB-like histones, respectively), increases the stability of recombinant
31 histone HFoB from *Methanobacterium formicicum* (Li et al. 2000). In addition, evidence
32 from mass spectrometry indicates that K35 in HTkB from *T. kodakarensis* and *Thermococcus*
33 *gammatolerans* is acetylated *in vivo* (Alpha-Bazin et al. 2021) although stoichiometry and
34 functional significance of this modification remain to be determined. A tyrosine at the same
35 position in HTkA removes the potential for acetylation. E18 in HMfB forms an
36 intermolecular salt bridge with K53, which helps to stabilise the interaction between
37 monomers in the histone dimer (Sandman et al. 2001; Decanniere et al. 2000). Mutating E18
38 to proline does not alter DNA binding (Soares et al. 2000) but loss of the intermonomer salt
39 bridge may result in less rigid dimer structures. Finally, having leucine (L) or phenylalanine
40 (F) at residue 46 has no obvious effect on DNA binding in HMfA/B, but the residue, located

1 at the interface between dimers, is important for tetramer formation (Soares et al. 2000; Marc
2 et al. 2002).

3
4 We considered how amino acid differences between HTkA- and HTkB-like histones affect
5 two key aspects of the histone-DNA complex: DNA affinity and tetramerization strength, a
6 proxy for tetramer stability. Using a structural modelling approach, we find that predicted
7 DNA binding for HTkB-like paralogs is, in most cases, stronger than for HTkA-like paralogs
8 (Fig 2, see Methods). This is in line with experimental findings that HTkB binds to DNA
9 more tightly and forms a more compact complex with DNA than HTkA (Higashibata et al.
10 1999). In contrast, predicted tetramerization strength does not strongly discriminate HTkA-
11 from HTkB-like histones (Fig 2). We also considered residues (K14, G17, K26, E30, E34,
12 Q48, E58, K61, K65) that were previously suggested to be important for stacking interactions
13 for either HTkA or HTkB, in the context of longer oligomers (Mattioli et al. 2017;
14 Henneman et al. 2018, 2021). None of these residues differ substantially between HTkA- and
15 HTkB-like histones, with the exception of Q48 (HTkB) (Henneman et al. 2018). We note that
16 Q48 is relatively uncommon in HTkB-like histones, and that *T. kodakarensis* HTkB may
17 therefore form more stable oligomeric complexes than the HTkB-like histones of most of its
18 cousins.

19 *HtkA- and HTkB-like histones form ancient paralogous groups*

20
21
22 Almost all Thermococcales have both an HTkA-like and an HTkB-like histone (Fig 1c). This
23 is consistent with (but not sufficient to demonstrate) ancient paralogy. To unravel the
24 evolutionary history of Thermococcales histones, we used RaxML-NG (Kozlov et al. 2019)
25 to build phylogenetic trees of all 123 Hmf-like histones found across the 61 genomes in our
26 analysis (see Methods). We find that HTkA-like and HTkB-like paralogs neatly separate into
27 two groups defined by their position on the tree (Fig 1d). This pattern of separation indicates
28 that one HTkA- and one HTkB-like histone were present in the last common ancestor of
29 Thermococcales. We detect only a small number of lineage-specific duplications and
30 deletions and find no evidence of rampant horizontal gene transfer (Fig 1c). The observation

1 that both paralogs have been maintained along divergent Thermococcales lineages strongly
2 suggests that at least some of the amino acid differences between them are functionally
3 important and under selection. Along with our recent report of ancient histone paralogs in the
4 Methanobacteriales (Stevens et al. 2020), this finding provides further evidence that histone
5 variants exist in archaea, evolving in parallel to those in eukaryotes. Note that, at present, we
6 have no convincing evidence that histone paralogs in the Thermococcales and those found in
7 the Methanobacteriales arose from the same ancient duplication event.

8 9 *Some Thermococcales encode histone-fold proteins similar to those found in bacteria*

10
11 Our survey also revealed that, alongside the HTkA/B-like histones, many Thermococcales
12 genomes encode histone-fold proteins similar to those found in some bacteria (Fig 1c, Table
13 S1), which harbour either a single or two (pseudodimeric) histone fold domains (Alva &
14 Lupas 2019). We will refer to these as bacteria-type singlets and doublets, respectively. Both
15 types are, on average, less well conserved than HTkA/B-like histones (Fig 3a) and their
16 distribution across the Thermococcales is noticeably patchier (Fig 1c, Fig 3f/g). Neither type
17 is present in the closest sister clades (Methanofastidiosa, Theionarchaea).

18
19 The bacteria-type singlet is confined to a monophyletic group that comprises some (but not
20 all) *Thermococcus* spp (Fig 1c, Fig 3f). It is present in all these species in the same highly
21 conserved syntenic context, suggesting a single origin (Fig 3f). Following acquisition, the
22 histone has been maintained along almost all lineages. There is only a single loss event (a
23 clean deletion) in the branch leading to *Thermococcus piezophilus* (Fig 3f). Conserved
24 synteny is also consistent with a single, earlier origin for the bacteria-type doublet, with
25 multiple subsequent losses (Fig 1c, Fig 3g). Based on the branching patterns and different
26 syntenic context, the bacteria-type histones in *Pyrococcus furiosus* have likely been acquired
27 secondarily from an extant or ancient *Thermococcus* species.

1 Bacteria-type doublets differ considerably in sequence from doublet histones previously
2 described in haloarchaea and *Methanopyrus kandleri* (Fig 3e). Outside the Thermococcales,
3 we only detect additional bacteria-type doublets in some (hyper)thermophilic
4 Methanocaldococcus and Archaeoglobus species, but never their mesophilic relatives,
5 suggestive of horizontal gene transfer in the high-temperature niche.

6
7 Bacteria-type histone-fold proteins have only recently been recognized and await functional
8 characterization. The only functional data we have at present comes from
9 transcriptome/proteome profiling. Consistent with lower sequence-level conservation (Fig
10 3a), the relative expression levels of these genes in *T. kodakarensis* (singlet: TK1040;
11 doublet: TK0750) are lower than those of HTkA/B-like histones at both the transcript and
12 protein level (Fig 3b) (Sas-Chen et al. 2020; Jäger et al. 2014). Together, they make up
13 0.37% of the measured exponential-phase proteome compared to 0.66% for HTkB and 1.1%
14 for HTkA. TK1040 was previously identified in *T. kodakarensis* chromatin fractions
15 (Maruyama et al. 2011), suggesting a (direct or indirect) association with DNA. However, the
16 same study estimated that less than 1% of the amount of chromatin-associated proteins were
17 attributable to TK1040. We therefore consider it unlikely that these histones are global
18 organizers of DNA similar to HTkA/B-like histones, but might modulate chromatin state,
19 either locally or globally, in response to environmental change. In both *P. furiosus* and *T.*
20 *kodakarensis*, the bacterial-type doublets are under the control of the heat shock regulator Phr
21 (encoded by PF1790 and TK2291, respectively) and upregulated upon Phr deletion,
22 suggesting a potential role in response to heat shock in these archaea (Kanai et al. 2010;
23 Keese et al. 2010). The *T. kodakarensis* doublet is also downregulated at lower temperatures,
24 similar to HTkA/B (Hocher et al. 2021), further consistent with a role in temperature
25 adaptation.

26
27 Can these histone fold proteins interact with HTkA/B-like histones? We used AlphaFold
28 (Jumper et al. 2021) to predict the structure of combinations of HTkA, HTkB and the
29 bacterial HF singlet from *T. kodakarensis*. Using this approach, all three are predicted to form
30 homodimers and, as expected, HTkA and HTkB form a stable heterodimer. When a
31 combination of either HTkA or HTkB and the bacterial singlet are used, Alphafold also
32 predicts that these will form a heterodimer (Fig 3c). The presence of non-HMf-like HF
33 proteins in Thermococcales genomes adds to the potential functional diversity of histone-
34 based chromatin in these species and may dynamically alter DNA accessibility at different

1 stages of cell growth or in response to environmental challenges. Further experimental
2 investigation is required, however, before meaningful conclusions can be drawn in this
3 regard, including whether they do interact, both structurally and functionally, with the
4 HTkA/B-like histones.

5 6 7 **Methods**

8 9 *Identification of histones in Thermococcales genomes*

10
11 Protein sets, genomes and GFF files for all available genomes of class Thermococci were
12 downloaded from GenBank (<https://www.ncbi.nlm.nih.gov/assembly>) using taxid 183968
13 [accessed on 2021-05-27]. Genomes not present in the GTDB tree (Parks et al. 2021)
14 (<https://gtdb.ecogenomic.org>, accessed 2021-08-01, see below) were removed. Two species
15 which were annotated as Thermococci in NCBI but branched outside the main group on the
16 GTDB tree were removed from the analysis, leaving a final set of 61 genomes, all from the
17 order Thermococcales. Protein sequences were predicted using Prodigal v2.6.3 (Hyatt et al.
18 2010) where not provided through GenBank. Histone proteins were extracted from the
19 protein sets through HMM searches using HMMER v3.3.1 (hmmsearch --noali) (Eddy 2011;
20 Finn et al. 2011) using Pfam models Cbfd_nfyb_hmf and DUF1931 (Finn et al. 2014) as
21 well as a Jackhmmer searches using the singlet and doublet histones from bacteria as a seed
22 used by others (Alva & Lupas 2019). Bacteria-type histones hit in the initial search were used
23 to build an HMM model and the Thermococcales protein set was re-searched using HMMER
24 v3.3.1 as above (hmmsearch --noali). Some proteins incorrectly identified as histones at this
25 stage were manually filtered out.

26 27 *Classification of Thermococcales histones into HtkA-like and HtkB-like groups*

28
29 Hmf-like histones in Thermococcales downloaded from GenBank [accessed on 2021-05-27]
30 (see above) were aligned using MAFFT (Katoh & Standley 2013) (-localpair --maxiterate
31 1000). Histones were clustered based on amino acid composition of their peptide sequences
32 using AASTats from the R package Seqinr (Charif & Lobry 2007). Histones which clustered
33 with HTkA and HTkB were assigned HTkA- or HTkB-like status, respectively (see Fig 1a).
34 20 maximum likelihood phylogenetic trees were built using Raxml-NG (Kozlov et al. 2019)

1 with the LG+G4 model of evolution as suggested by ModelTest-NG (Darriba et al. 2020).
2 The unrooted best maximum likelihood tree is shown. All trees were plotted using iTOL
3 (Letunic & Bork 2019). Orthologous genes in the genomic neighbourhood of each histone
4 were highlighted on the tree using Genespy as best reciprocal hits (Garcia et al. 2019).
5 Orthologs in the histone neighbourhood were identified by performing reciprocal best hits for
6 each genome against *T. kodakarensis* using BLAST (Altschul et al. 1990), retaining those
7 which have a similarity score of >40% and are within 20% length of one another (Rocha
8 2006). Note that the reciprocal best hit approach was only applied to identify putative
9 orthologs in the neighbourhood of histone genes. The histones themselves were identified
10 using HMMer searches (see above) and subsequent analyses not restricted to reciprocal best
11 hits of HTkA/B. To generate sequence logos, histones were aligned using MAFFT (Katoh &
12 Standley 2013) (`-localpair --maxiterate 1000`) and visualized using ggseqlogo in R (Wagih
13 2017).

15 *Predicted DNA binding and tetramerisation*

17 Predicted DNA-binding and interaction (tetramerization) strength between dimers for
18 Thermococcales species was computed as in (Stevens et al. 2020). In brief, sequences were
19 aligned to HMfB and substitutions were mapped onto a tetrameric model of HMfB (extracted
20 from PDB structure 5t5k) using FoldX (Schymkowitz et al. 2005) to generate models of
21 homotetramers with DNA. Structures were then energy-minimised using AmberTools (Maier
22 et al. 2015) and binding affinity was calculated using an MMPBSA approach with the ff14SB
23 forcefield (Miller et al. 2012). $\Delta\Delta G$ was calculated relative to HMfB. The mean value for five
24 replicates is shown for each model.

26 *Species tree*

28 The archaeal species tree was downloaded from GTDB (<https://gtdb.ecogenomic.org>) on
29 2021-08-01.

31 *Expression data*

1 Expression data for *T. kodakarensis* was obtained from primary sources and NCBI's Gene
2 Expression Omnibus (GEO) (Barrett et al. 2013). Proteomics data from (Sas-Chen et al.
3 2020) was processed as in Hocher et al. (2021). Protein abundance at 85°C is shown.
4 Transcript abundance data was obtained from (Jäger et al. 2014) and is shown as transcripts
5 per million (TPM).

6 *Bacteria-type singlet histone fold and HTkA/B structure prediction*

7
8
9 The AlphaFold v2.0 (Jumper et al. 2021) collab notebook
10 ([https://colab.research.google.com/github/deepmind/alphafold/blob/main/notebooks/AlphaFo](https://colab.research.google.com/github/deepmind/alphafold/blob/main/notebooks/AlphaFold.ipynb)
11 [ld.ipynb](https://colab.research.google.com/github/deepmind/alphafold/blob/main/notebooks/AlphaFold.ipynb) accessed 18-08-21) was used to predict the structure of HTkA, HTkB and the
12 bacterial singlet HF protein (TK1040) as homodimers and all heterodimer combinations. The
13 MSA method used was jackhmmer, and models were ranked by PTMscore. The top model is
14 shown for all homodimers and heterodimers. Images shown were generated using UCSF,
15 ChimeraX developed by the Resource for Biocomputing, Visualization, and Informatics at
16 the University of California, San Francisco, with support from National Institutes of Health
17 (R01-GM129325) and the Office of Cyber Infrastructure and Computational Biology,
18 National Institute of Allergy and Infectious Diseases (Pettersen et al. 2021).

19 20 *Identification of doublet histones in archaea*

21
22 The genomes, protein sets and GFF files for a balanced set of archaea species was
23 downloaded on 2021-05-21 from GenBank (<https://www.ncbi.nlm.nih.gov/assembly>) using
24 taxid 2157, and processed as above (see: *Identification of histones in Thermococcales*
25 *genomes*). Doublet histones, including doublets from Halobacteria (Dulmage et al. 2015) and
26 HMK in *Methanopyrus kandleri* (Fahrner et al. 2001; Slesarev et al. 1998) were identified by
27 their position on a CDS-level tree of archaeal histones, and by length. This larger tree was
28 built from 20 ML trees using Raxml-NG (Kozlov et al. 2019). Trees for the doublet histones
29 were then built as previously described (see: *Classification of Thermococcales histones into*
30 *HtkA-like and HtkB-like groups*) and all figures plotted using iTOL (Letunic & Bork 2019).
31 Annotation for the orthologous genes in the genomic neighbourhood was generated using

1 GeneSpy (Garcia et al. 2019). Clades with doublet histones are annotated on a tree of archaea
2 adapted from (Borrel et al. 2020) to include Hodarchaeota (Liu et al. 2021).

3 4 *Percentage identity of HTkA-/HTkB-like and bacteria-type histones*

5
6 HTkA-like, HTkB-like and bacteria-type histones were aligned using MAFFT (Katoh &
7 Standley 2013) (`--localpair --maxiterate 1000`) separately for species containing either
8 bacteria-type doublet or bacteria-type singlet histones. For each histone type, pairwise
9 sequence identity was calculated using `seqidentity` from the R package `bio3d` (Grant et al.
10 2006).

11 12 **Data Availability Statement**

13
14 All data underlying were derived from sources in the public domain: GenBank
15 (<https://www.ncbi.nlm.nih.gov/assembly>); GTDB (<https://gtdb.ecogenomic.org>); and data
16 from original publications cited throughout the manuscript.

17 18 **References**

- 19
20 Alpha-Bazin B et al. 2021. Lysine-specific acetylated proteome from the archaeon
21 *Thermococcus gammatolerans* reveals the presence of acetylated histones. *J. Proteomics*.
22 232:104044. doi: 10.1016/j.jprot.2020.104044.
- 23 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search
24 tool. *J. Mol. Biol.* 215:403–410. doi: 10.1016/S0022-2836(05)80360-2.
- 25 Alva V, Lupas AN. 2019. Histones Predate the Split Between Bacteria and Archaea.
26 *Bioinformatics*. 35:2349–2353. doi: 10.1093/bioinformatics/bty1000.
- 27 Bailey KA, Marc F, Sandman K, Reeve JN. 2002. Both DNA and histone fold sequences
28 contribute to archaeal nucleosome stability. *J. Biol. Chem.* 277:9293–9301. doi:
29 10.1074/jbc.M110029200.
- 30 Bailey KA, Pereira SL, Widom J, Reeve JN. 2000. Archaeal histone selection of nucleosome
31 positioning sequences and the procaryotic origin of histone-dependent genome evolution. *J.*
32 *Mol. Biol.* 303:25–34. doi: 10.1006/jmbi.2000.4128.
- 33 Barrett T et al. 2013. NCBI GEO: Archive for functional genomics data sets - Update.

- 1 Nucleic Acids Res. 41:991–995. doi: 10.1093/nar/gks1193.
- 2 Bönisch C, Hake SB. 2012. Histone H2A variants in nucleosomes and chromatin: More or
3 less stable? Nucleic Acids Res. 40:10719–10741. doi: 10.1093/nar/gks865.
- 4 Borrel G, Brugère JF, Gribaldo S, Schmitz RA, Moissl-Eichinger C. 2020. The host-
5 associated archaeome. Nat. Rev. Microbiol. doi: 10.1038/s41579-020-0407-y.
- 6 Bowerman S, Wereszczynski J, Luger K. 2021. Archaeal chromatin ‘slinkies’ are inherently
7 dynamic complexes with deflected DNA wrapping pathways. Elife. 10:e65587. doi:
8 10.7554/eLife.65587.
- 9 Charif D, Lobry JR. 2007. SeqinR 1.0-2: a contributed package to the R project for statistical
10 computing devoted to biological sequences retrieval and analysis. In: Structural approaches to
11 sequence evolution: Molecules, networks, populations. Bastolla, U, Porto, M, Roman, HE, &
12 Vendruscolo, M, editors. Springer Verlag: New York pp. 207–232.
- 13 Čuboňová L et al. 2012. An archaeal histone is required for transformation of *Thermococcus*
14 *kodakarensis*. J. Bacteriol. 194:6864–6874. doi: 10.1128/JB.01523-12.
- 15 Darriba Di et al. 2020. ModelTest-NG: A New and Scalable Tool for the Selection of DNA
16 and Protein Evolutionary Models. Mol. Biol. Evol. 37:291–294. doi:
17 10.1093/molbev/msz189.
- 18 Decanniere K, Babu AM, Sandman K, Reeve JN, Heinemann U. 2000. Crystal structures of
19 recombinant histones HMfA and HMfB from the hyperthermophilic archaeon
20 *Methanothermus fervidus*. J. Mol. Biol. 303:35–47. doi: 10.1006/jmbi.2000.4104.
- 21 Dulmage KA, Todor H, Schmid AK. 2015. Growth-Phase-Specific Modulation of Cell
22 Morphology and Gene Expression by an Archaeal Histone Protein. MBio. 6. doi:
23 10.1128/mbio.00649-15.
- 24 Eddy SR. 2011. Accelerated Profile HMM Searches. PLoS Comput. Biol. 7:e1002195. doi:
25 10.1371/journal.pcbi.1002195.
- 26 Fahrner RL, Cascio D, Lake JA, Slesarev A. 2001. An ancestral nuclear protein assembly:
27 Crystal structure of the *Methanopyrus kandleri* histone. Protein Sci. 10:2002–2007. doi:
28 10.1110/ps.10901.
- 29 Finn RD et al. 2014. Pfam: The protein families database. Nucleic Acids Res. 42:222–230.
30 doi: 10.1093/nar/gkt1223.
- 31 Finn RD, Clements J, Eddy SR. 2011. HMMER web server: Interactive sequence similarity
32 searching. Nucleic Acids Res. 39:29–37. doi: 10.1093/nar/gkr367.
- 33 Garcia PS, Jauffrit F, Grangeasse C, Brochier-Armanet C. 2019. GeneSpy, a user-friendly
34 and flexible genomic context visualizer. Bioinformatics. 35:329–331. doi:

- 1 10.1093/bioinformatics/bty459.
- 2 Grant BJ, Rodrigues APC, ElSawy KM, McCammon JA, Caves LSD. 2006. Bio3d: an R
- 3 package for the comparative analysis of protein structures. *Bioinformatics*. 22:2695–2696.
- 4 doi: 10.1093/bioinformatics/btl461.
- 5 Henikoff S, Smith MM. 2015. Histone Variants and Epigenetics. *Cold Spring Harb. Perspect.*
- 6 *Biol.* 7:a019364. doi: 10.1101/cshperspect.a019364.
- 7 Henneman B. 2019. Histone-DNA assemblies in archaea: shaping the genome on the edge of
- 8 life. Leiden University.
- 9 Henneman B et al. 2021. Mechanical and structural properties of archaeal hypernucleosomes.
- 10 *Nucleic Acids Res.* 49:4338–4349. doi: 10.1093/nar/gkaa1196.
- 11 Henneman B, van Emmerik C, van Ingen H, Dame RT. 2018. Structure and function of
- 12 archaeal histones. *PLoS Genet.* 14:e1007582. doi: 10.1371/journal.pgen.1007582.
- 13 Higashibata H, Fujiwara S, Takagi M, Imanaka T. 1999. Analysis of DNA compaction
- 14 profile and intracellular contents of archaeal histones from *Pyrococcus kodakaraensis* KOD1.
- 15 *Biochem. Biophys. Res. Commun.* 258:416–424. doi: 10.1006/bbrc.1999.0533.
- 16 Hocher A et al. 2021. Growth temperature is the principal driver of chromatinization in
- 17 archaea. *BioRxiv*. doi: <https://doi.org/10.1101/2021.07.08.451601>.
- 18 Hyatt D et al. 2010. Prodigal: Prokaryotic gene recognition and translation initiation site
- 19 identification. *BMC Bioinformatics*. 11:119. doi: 10.1186/1471-2105-11-119.
- 20 Jäger D, Förstner KU, Sharma CM, Santangelo TJ, Reeve JN. 2014. Primary transcriptome
- 21 map of the hyperthermophilic archaeon *Thermococcus kodakarensis*. *BMC Genomics*.
- 22 15:684. doi: 10.1186/1471-2164-15-684.
- 23 Jang C, Shibata Y, Starmer J, Yee D, Magnuson T. 2015. Histone H3.3 maintains genome
- 24 integrity during mammalian development. *Genes Dev.* 29:1377–1392. doi:
- 25 10.1101/gad.264150.115.GENES.
- 26 Jumper J et al. 2021. Highly accurate protein structure prediction with AlphaFold. *Nature*.
- 27 596:583–589. doi: 10.1038/s41586-021-03819-2.
- 28 Kanai T, Takedomi S, Fujiwara S, Atomi H, Imanaka T. 2010. Identification of the Phr-
- 29 dependent heat shock regulon in the hyperthermophilic archaeon, *Thermococcus*
- 30 *kodakaraensis*. *J. Biochem.* 147:361–370. doi: 10.1093/jb/mvp177.
- 31 Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7:
- 32 Improvements in performance and usability. *Mol. Biol. Evol.* 30:772–780. doi:
- 33 10.1093/molbev/mst010.
- 34 Keese AM, Schut GJ, Ouhammouch M, Adams MWW, Thomm M. 2010. Genome-wide

- 1
2
3 1 identification of targets for the archaeal heat shock regulator Phr by cell-free transcription of
4 genomic DNA. *J. Bacteriol.* 192:1292–1298. doi: 10.1128/JB.00924-09.
5
6 3 Kozlov AM, Darriba D, Flouri T, Morel B, Stamatakis A. 2019. RAxML-NG: A fast,
7
8 4 scalable and user-friendly tool for maximum likelihood phylogenetic inference.
9
10 5 *Bioinformatics.* 35:4453–4455. doi: 10.1093/bioinformatics/btz305.
11
12 6 Letunic I, Bork P. 2019. Interactive Tree of Life (iTOL) v4: Recent updates and new
13
14 7 developments. *Nucleic Acids Res.* 47:256–259. doi: 10.1093/nar/gkz239.
15
16 8 Li WT, Shriver JW, Reeve JN. 2000. Mutational analysis of differences in thermostability
17
18 9 between histones from mesophilic and hyperthermophilic archaea. *J. Bacteriol.* 182:812–817.
19
20 10 doi: 10.1128/JB.182.3.812-817.2000.
21
22 11 Liu Y et al. 2021. Expanded diversity of Asgard archaea and their relationships with
23
24 12 eukaryotes. *Nature.* 593:553–557. doi: 10.1038/s41586-021-03494-3.
25
26 13 Maier JA et al. 2015. ff14SB: Improving the Accuracy of Protein Side Chain and Backbone
27
28 14 Parameters from ff99SB. *J. Chem. Theory Comput.* 11:3696–3713. doi:
29
30 15 10.1021/acs.jctc.5b00255.
31
32 16 Marc F, Sandman K, Lurz R, Reeve JN. 2002. Archaeal histone tetramerization determines
33
34 17 DNA affinity and the direction of DNA supercoiling. *J. Biol. Chem.* 277:30879–30886. doi:
35
36 18 10.1074/jbc.M203674200.
37
38 19 Martire S, Banaszynski LA. 2020. The roles of histone variants in fine-tuning chromatin
39
40 20 organization and function. *Nat. Rev. Mol. Cell Biol.* 21:522–541. doi: 10.1038/s41580-020-
41
42 21 0262-8.
43
44 22 Maruyama H et al. 2013. An alternative beads-on-a-string chromatin architecture in
45
46 23 *Thermococcus kodakarensis*. *EMBO Rep.* 14:711–717. doi: 10.1038/embor.2013.94.
47
48 24 Maruyama H et al. 2011. Histone and TK0471/TrmBL2 form a novel heterogeneous genome
49
50 25 architecture in the hyperthermophilic archaeon *Thermococcus kodakarensis*. *Mol. Biol. Cell.*
51
52 26 22:386–398. doi: 10.1091/mbc.e10-08-0668.
53
54 27 Mattioli F et al. 2017. Structure of histone-based chromatin in Archaea. *Science.* 357:609–
55
56 28 612. doi: 10.1126/science.aaj1849.
57
58 29 Miller BR et al. 2012. MMPBSA.py: An efficient program for end-state free energy
59
60 30 calculations. *J. Chem. Theory Comput.* 8:3314–3321. doi: 10.1021/ct300418h.
31
32 31 Nalabothula N et al. 2013. Archaeal nucleosome positioning in vivo and in vitro is directed
33
34 32 by primary sequence motifs. *BMC Genomics.* 14:391. doi: 10.1186/1471-2164-14-391.
35
36 33 Palmer DK, O’Day K, Trong HLE, Charbonneau H, Margolis RL. 1991. Purification of the
37
38 34 centromere-specific protein CENP-A and demonstration that it is a distinctive histone. *Proc.*

- 1 Natl. Acad. Sci. U. S. A. 88:3734–3738. doi: 10.1073/pnas.88.9.3734.
- 2 Parks DH et al. 2021. GTDB: an ongoing census of bacterial and archaeal diversity through a
3 phylogenetically consistent, rank normalized and complete genome-based taxonomy. *Nucleic*
4 *Acids Res.* doi: 10.1093/nar/gkab776.
- 5 Pettersen EF et al. 2021. UCSF ChimeraX: Structure visualization for researchers, educators,
6 and developers. *Protein Sci.* 30:70–82. doi: 10.1002/pro.3943.
- 7 Rocha EPC. 2006. Inference and analysis of the relative stability of bacterial chromosomes.
8 *Mol. Biol. Evol.* 23:513–522. doi: 10.1093/molbev/msj052.
- 9 Rojcek M, Hocher A, Stevens KM, Merckenschlager M, Warnecke T. 2019. Chromatinization
10 of *Escherichia coli* with archaeal histones. *Elife.* 8:e49038. doi: 10.7554/eLife.49038.
- 11 Sanders TJ et al. 2021. Extended Archaeal Histone-Based Chromatin Structure Regulates
12 Global Gene Expression in *Thermococcus kodakarensis*. *Front. Microbiol.* 12:1071. doi:
13 10.3389/fmicb.2021.681150.
- 14 Sanders TJ et al. 2019. TFS and Spt4/5 accelerate transcription through archaeal histone-
15 based chromatin. *Mol. Microbiol.* 111:784–797. doi: 10.1111/mmi.14191.
- 16 Sandman K, Grayling RA, Dobrinski B, Lurz R, Reeve JN. 1994. Growth-phase-dependent
17 synthesis of histones in the archaeon *Methanothermus fervidus*. *Proc. Natl. Acad. Sci.*
18 91:12624–12628. doi: 10.1073/pnas.91.26.12624.
- 19 Sandman K, Krzycki JA, Dobrinski B, Lurz R, Reeve JN. 1990. HMf, a DNA-binding
20 protein isolated from the hyperthermophilic archaeon *Methanothermus fervidus*, is most
21 closely related to histones. *Proc. Natl. Acad. Sci.* 87:5788–5791. doi:
22 10.1073/pnas.87.15.5788.
- 23 Sandman K, Soares D, Reeve JN. 2001. Molecular components of the archaeal nucleosome.
24 *Biochimie.* 83:277–281. doi: 10.1016/S0300-9084(00)01208-6.
- 25 Sas-Chen A et al. 2020. Dynamic RNA acetylation revealed by quantitative cross-
26 evolutionary mapping. *Nature.* 583:638–643. doi: 10.1038/s41586-020-2418-2.
- 27 Schymkowitz J et al. 2005. The FoldX web server: An online force field. *Nucleic Acids Res.*
28 33:W382–W388. doi: 10.1093/nar/gki387.
- 29 Sitbon D, Boyarchuk E, Dingli F, Loew D, Almouzni G. 2020. Histone variant H3.3 residue
30 S31 is essential for *Xenopus* gastrulation regardless of the deposition pathway. *Nat.*
31 *Commun.* 11:1256. doi: 10.1038/s41467-020-15084-4.
- 32 Slesarev AI, Belova GI, Kozyavkin SA, Lake JA. 1998. Evidence for an early prokaryotic
33 origin of histones H2A and H4 prior to the emergence of eukaryotes. *Nucleic Acids Res.*
34 26:427–430. doi: 10.1093/nar/26.2.427.

- 1 Soares DJ, Sandman K, Reeve JN. 2000. Mutational analysis of archaeal histone-DNA
 2 interactions. *J. Mol. Biol.* 297:39–47. doi: 10.1006/jmbi.2000.3546.
- 3 Starich MR, Sandman K, Reeve JN, Summers MF. 1996. NMR structure of HMfB from the
 4 hyperthermophile, *Methanothermus fervidus*, confirms that this archaeal protein is a histone.
 5 *J. Mol. Biol.* 255:187–203. doi: 10.1006/jmbi.1996.0016.
- 6 Stevens KM et al. 2020. Histone variants in archaea and the evolution of combinatorial
 7 chromatin complexity. *Proc. Natl. Acad. Sci. U. S. A.* 117:33384–33395. doi:
 8 10.1073/PNAS.2007056117.
- 9 Talbert PB, Henikoff S. 2010. Histone variants - ancient wrap artists of the epigenome. *Nat.*
 10 *Rev. Mol. Cell Biol.* 11:264–275. doi: 10.1038/nrm2861.
- 11 Wagih O. 2017. ggseqlogo: A versatile R package for drawing sequence logos.
 12 *Bioinformatics.* 33:3645–3647. doi: 10.1093/bioinformatics/btx469.
- 13 Wolfe JM, Fournier GP. 2018. Horizontal gene transfer constrains the timing of methanogen
 14 evolution. *Nat. Ecol. Evol.* 2:897–903. doi: 10.1038/s41559-018-0513-7.

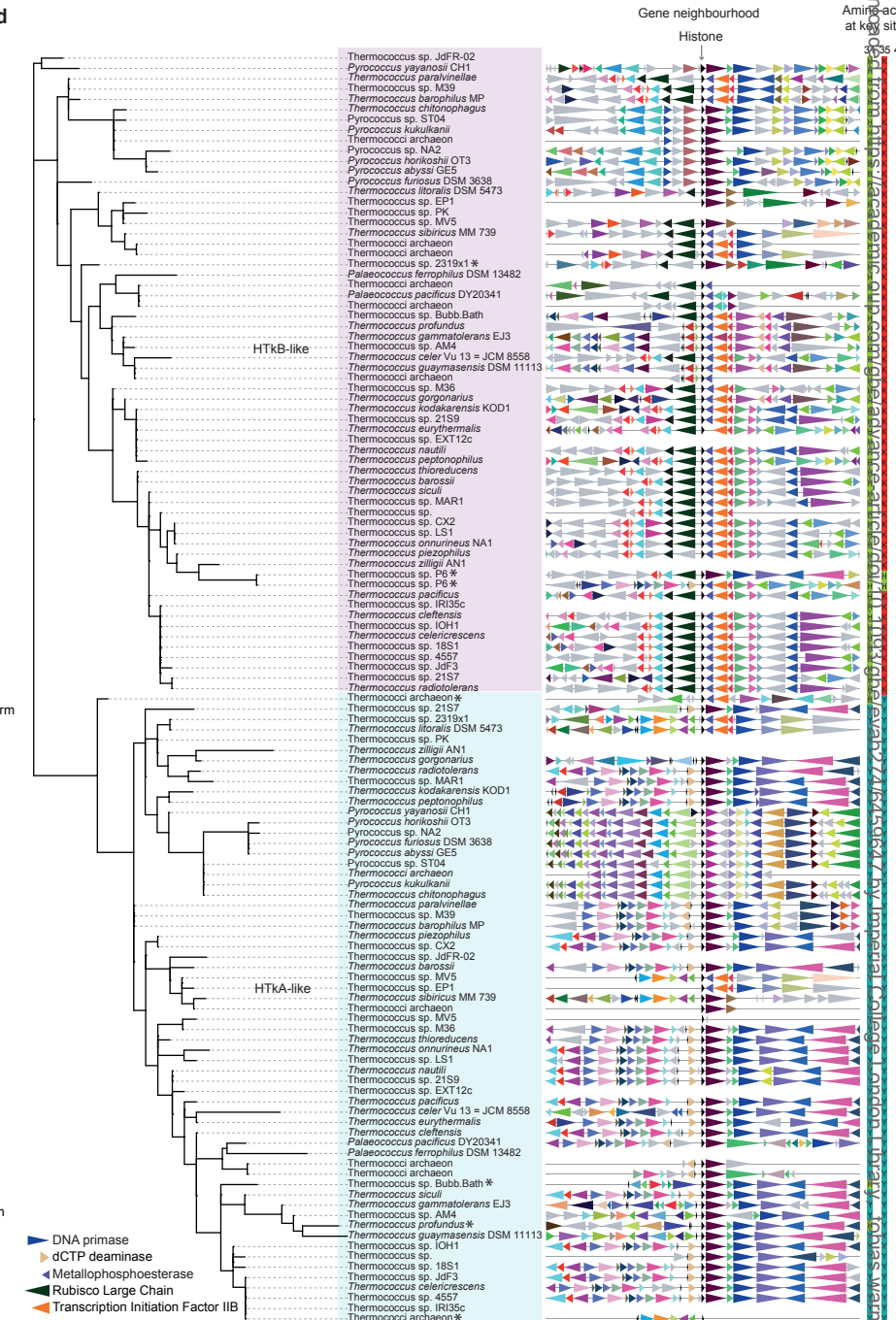
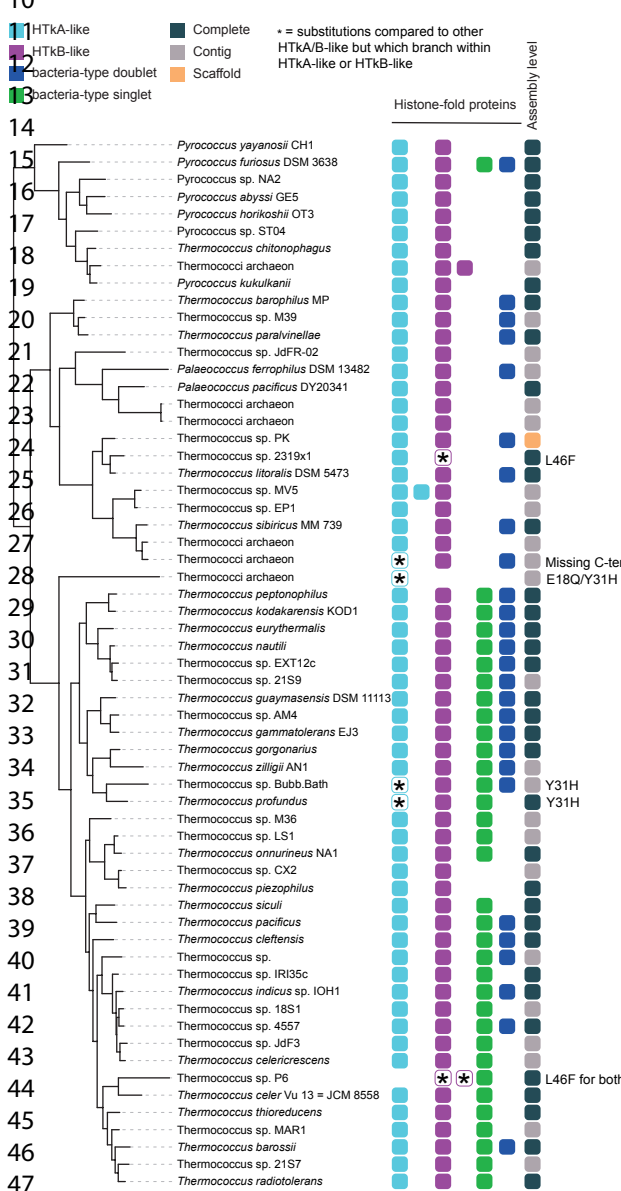
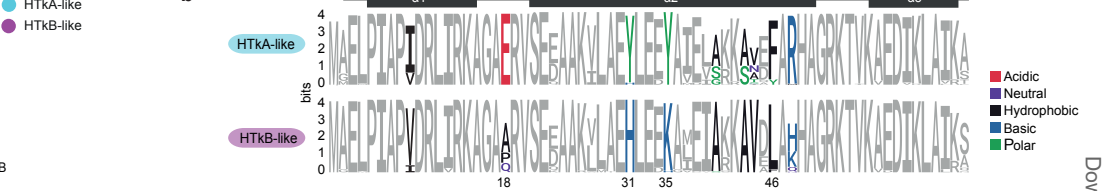
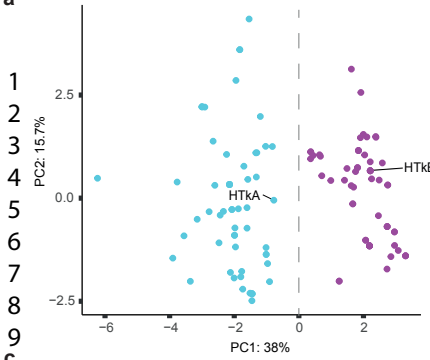
15 **Figure legends**

16
 17
 18 **Figure 1. a.** Principle component analysis of Thermococcales HMf-like histones based on
 19 amino acid properties (AAStats, see Methods). Histones that cluster with either HTkA or
 20 HTkB along the first principal component are coloured accordingly (HTkA-like histones in
 21 blue, HTkB-like in purple). **b.** Sequence logos showing amino acid composition of HTkA-
 22 and HTkB-like histones across 61 Thermococcales. Amino acids that differ substantially
 23 between the two groups are coloured. Positions are numbered relative to HMfB from
 24 *Methanothermus fervidus* to facilitate comparison with prior studies. **c.** GTDB species tree
 25 for all Thermococcales in the dataset indicating presence/absence of histones of a particular
 26 type in each genome. Each square represents one histone and is coloured by histone type. **d.**
 27 Protein-level phylogenetic tree of all HMf-like Thermococcales histones. Genes in the
 28 neighbourhood are coloured to indicate ortholog identity (see Methods). Note that, while the
 29 gene neighbourhood is broadly conserved for HTkA- vis-à-vis HTkB-like orthologs, some
 30 HTkB-like genes (e.g. in *Thermococcus chitonophagus*) have a 3' neighbourhood normally
 31 found for HTkA. This is owing to a large genomic rearrangement event in which HTkB
 32 served as the breakpoint (not shown).

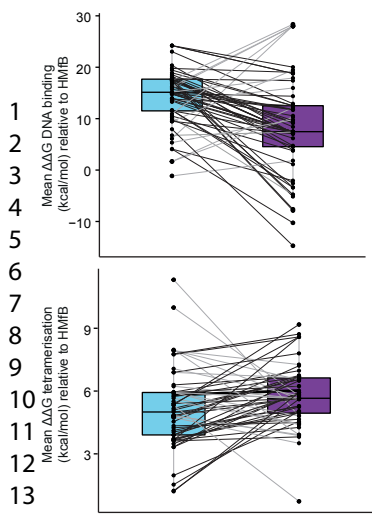
33
 34 **Figure 2.** Predicted DNA binding affinity (top) and tetramer stability (bottom) for HTkA/B-

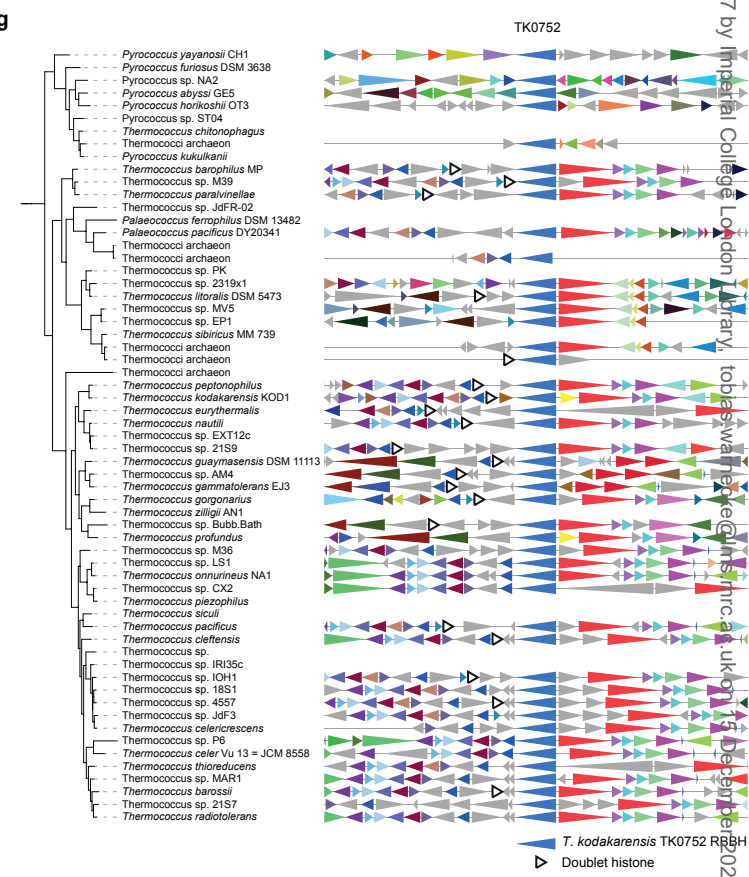
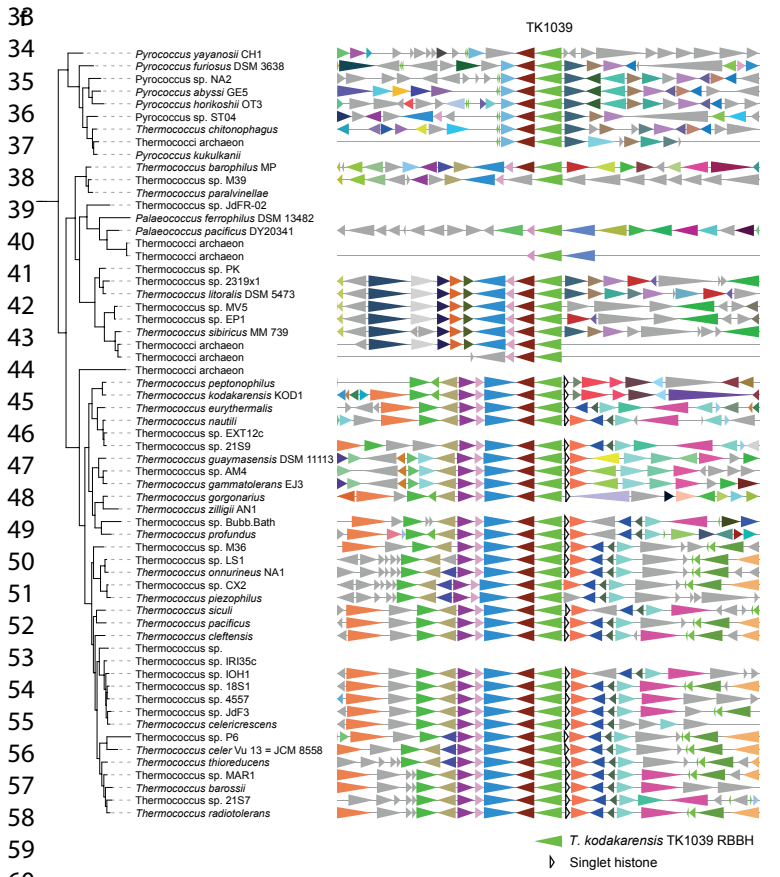
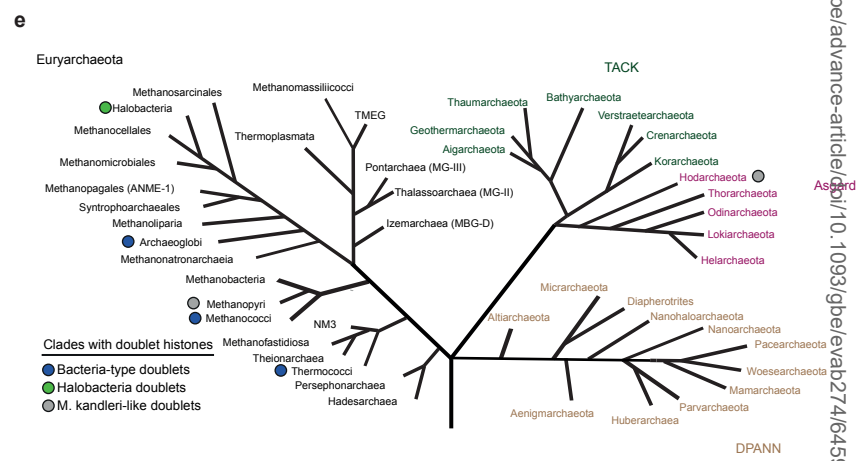
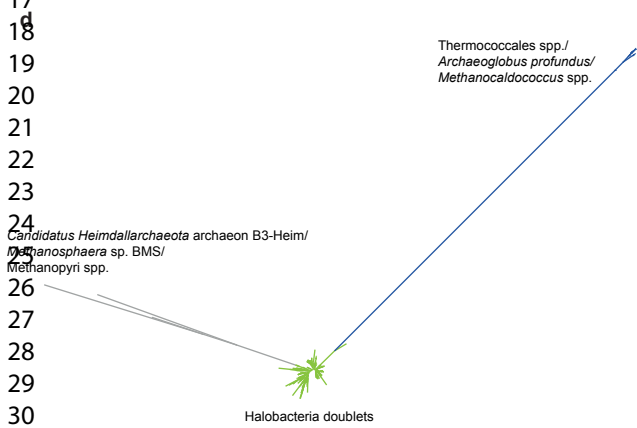
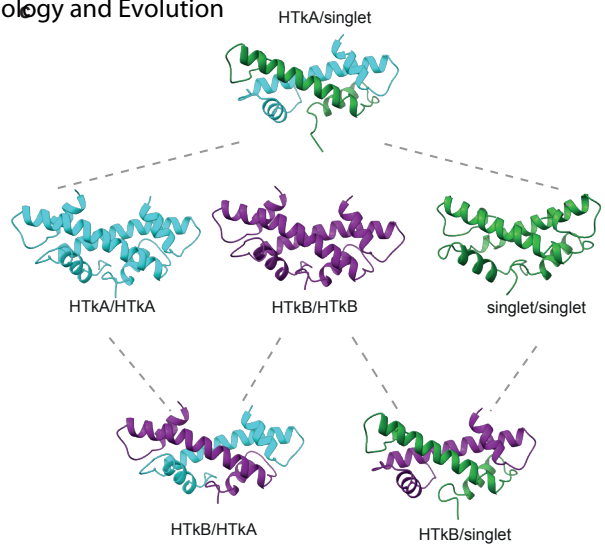
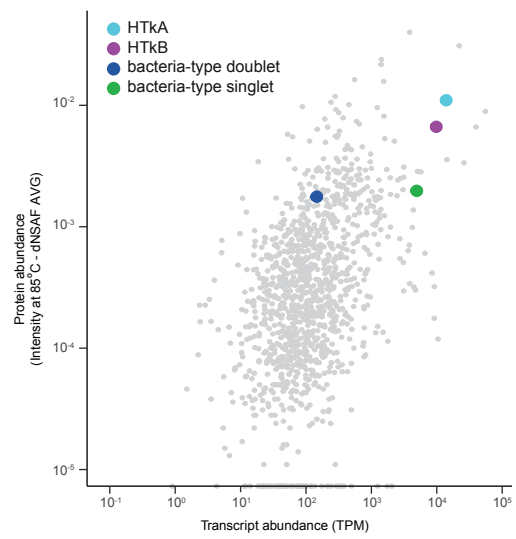
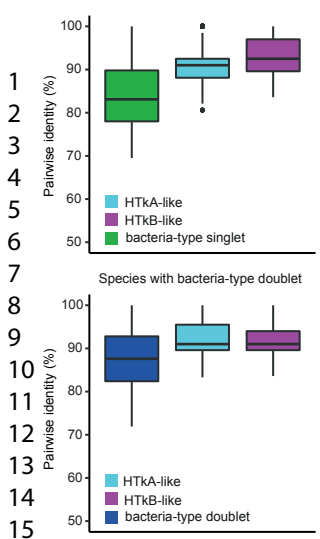
1 like paralogs. Lines connect paralogs from the same genome. Lines are black when DNA
2 binding affinity is stronger or tetramer interface energy weaker, respectively, for the HTkB-
3 like paralog.

4
5 **Figure 3. a.** Pairwise identity for histones in species with HTkA-like, HTkB-like, and
6 bacteria-type singlet (top) or doublet histones (bottom) across all Thermococcales. **b.** Protein
7 and transcript abundance (see Methods) of genes in *T. kodakarensis*. TPM: transcripts per
8 million. **c.** AlphaFold-predicted homodimeric structures of TK1040 (green), HTkA (light
9 blue) and HTkB (purple) and heterodimers of HTkA and HTkB (light blue/purple), HTkA
10 and the bacteria-type singlet TK1040 (light blue/green), and HTkB and TK1040
11 (purple/green). **d.** Protein-level tree for doublet histones (containing an end-to-end
12 duplication) in archaea. Green branches contain proteins from Halobacteria, blue branches
13 contain bacteria-type doublet histones in Thermococcales, *Methanocaldococcus* spp. and one
14 *Archaeoglobus* sp., grey branches show doublet histones including HMk from *Methanopyrus*
15 *kandleri* (Fahrner et al. 2001; Slesarev et al. 1998). **e.** Clades with doublet histones
16 highlighted on a tree capturing archaeal diversity, adapted from (Borrel et al. 2020; Stevens
17 et al. 2020) to include the clade for the species *Candidatus Heimdallarchaeota* archaeon B3-
18 Heim, more recently suggested to be a member the Hodarchaeota (Liu et al. 2021). Circles
19 denote clades containing a doublet histone, and are coloured by the type of doublet (see d).
20 **f,g.** Thermococcales GTDB species tree showing syntenic regions for TK1039 (light green, in
21 **f**) and TK0752 (blue, in **g**) commonly found located close to the bacteria-type singlet or
22 doublet histones. Genes are coloured to indicate homology (see Fig 1d/Methods).



48
49
50
51
52
53
54
55
56
57
58
59
60





Downloaded from https://academic.oup.com/gbe/advance-article/doi/10.1093/gbe/evab274/6459647 by Imperial College London Library, on 13 December 2021