

The histone fold: A ubiquitous architectural motif utilized in DNA compaction and protein dimerization

(paired-element motif/archaeal histones/centromeric CENP-A/transcription/evolution)

GINA ARENTS* AND EVANGELOS N. MOUDRIANAKIS*†

*Department of Biology, The Johns Hopkins University, Baltimore, MD 21218; and †Biology Department, University of Athens, Athens, Greece

Communicated by Hamilton O. Smith, Johns Hopkins School of Medicine, Baltimore, MD, August 21, 1995 (received for review July 6, 1995)

ABSTRACT The histones of all eukaryotes show only a low degree of primary structure homology, but our earlier crystallographic results defined a three-dimensional structural motif, the histone fold, common to all core histones. We now examine the specific architectural patterns within the fold and analyze the nature of the amino acid residues within its functional segments. The histone fold emerges as a fundamental protein dimerization motif while the differentiations of the tips of the histone dimers appear to provide the rules of core octamer assembly and the basis for nucleosome regulation. We present evidence for the occurrence of the fold from archaeobacteria to mammals and propose the use of this structural motif to define a distinct family of proteins, the histone fold superfamily. It appears that evolution has conserved the conformation of the fold even through variations in primary structure and among proteins with various functional roles.

The fundamental structural unit of all chromosomes, the nucleosome, comprises an almost constant length of DNA wrapped tightly around a protein spool, the core histone octamer. Earlier, we solved the crystal structure of the histone octamer (1) and have found it to be a tripartite assembly in which two (H2A–H2B) dimers flank a centrally located (H3–H4)₂ tetramer. The four types of core histone chains have very low sequence homology but share a common motif of tertiary structure, the *histone fold* (1). In this study we present insights into the organization of the histone fold, its persistence through all core histones from archaeobacteria to mammals, and its involvement in generating histone heterodimers via the *handshake motif* (1) of assembly. Furthermore, we discuss the architectural and evolutionary attributes of this motif in relation to its role in DNA condensation and gene regulation. Recently (2), through an extensive search of protein sequence data banks we identified a consensus primary structure for the histone fold and found it present, from bacteria to mammals, in a large number of proteins with diverse functions (e.g., transcription factors, enzymes, etc.). We propose that the attributes of the histone fold can be used as a ruler for defining a distinct protein superfamily.

METHODS

The histone fold is found interstitially within each histone chain and at different absolute locations from the starting residue of each chain. To facilitate comparison of equivalent amino acid residues within the four fold motifs, we have identified the residues within the fold using the prefix *F* followed by the amino acid identity symbol and the sequential number of its location within the fold beginning with *F1* for the first fold residue. Thus, *FLys-49* denotes a lysine residue at the 49th position within the fold region of a histone chain. The

histone fold regions of the four histone chains will be referred to as *FH2A*, *FH2B*, *FH3*, and *FH4*. All comparisons between histone structures were performed using the program ALIGN, kindly made available to us by Mario Amzel. ALIGN calculates the best fit and rms deviation between sets of equivalent α carbon coordinates. The histone coordinates employed in the calculation are taken from the 3.1-Å, partially refined ($R = 26\%$), histone octamer structure (1). In performing the calculations, allowance was made for a single-site deletion in the loop region of H4 that follows helix I. The H4 deletion has been arbitrarily assigned to residue 15 of the fold, since neither a comparison of the structures nor an examination of the amino acid sequences made obvious whether the deletion is, in fact, at residue 14 or 15 of the fold. Accordingly, we have ignored positions 14 and 15 during the alignment and comparisons of the four core histone structures.

RESULTS

Secondary Structure of the Histone Fold. The four classes of the core histones (H2A, H2B, H3, H4) contain three types of structural motifs: (i) the histone fold, (ii) the extra-fold structured elements unique to the different histones, and (iii) the labile termini, which vary in length from 13 to 42 amino acids (1). The histone fold consists of an 11-residue helix (helix I), followed by a short loop and β strand (strand A), a long 27-residue helix (helix II), another short loop and β strand (strand B), and a final 11-residue helix (helix III). At the current level of resolution (3.1 Å) and refinement ($R = 26\%$), the exact number of residues in each helix and loop or strand segment appears to vary by one or two from histone to histone. For comparison, the sequences of the fold regions of the four core histones are shown in Fig. 1. While only 4% identity exists when all four chains in the fold region are considered simultaneously, this increases by a factor of 4–5 when the fold regions are compared in all six pairwise combinations.

Three-Dimensional Structure of the Histone Fold. The three-dimensional folding patterns of the four core histones are remarkably similar, as shown in Fig. 2, more so than could be predicted from comparisons of their primary structures. We have compared analogous α carbon positions in the fold domain of each core histone after pairwise alignment with another fold in all possible permutations (Table 1). *FH2B*, *FH3*, and *FH4* exhibit the greatest similarity to one another, while *FH2A* differs the most from the other three histones, largely due to the somewhat altered orientation of helix I (Fig. 2). The variation of helix I in *FH2A* follows a frequently observed type of protein structure wobble about loop segments of a protein (6).

We reported earlier that each histone fold appears to be the result of a tandem duplication that divides it into two similar and contiguous helix–strand–helix (HSH) motifs (7). Here, we refer to the amino half of the histone fold as HSH1 and to the carboxyl half as HSH2, as we present a detailed comparison of

The publication costs of this article were defrayed in part by page charge payment. This article must therefore be hereby marked "advertisement" in accordance with 18 U.S.C. §1734 solely to indicate this fact.

Abbreviations: HSH, helix–strand–helix; PEM, paired-element motif.

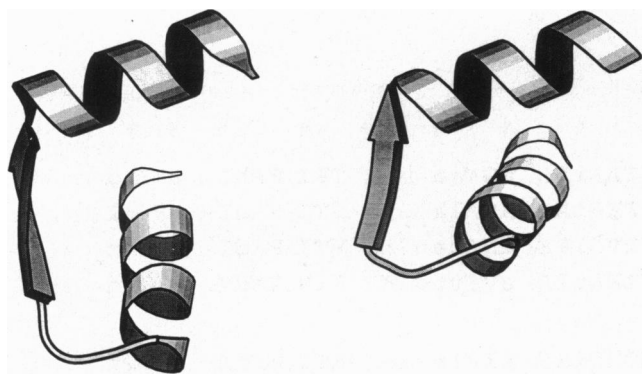


FIG. 3. Comparison of the amino and carboxyl halves of the histone fold. The example here is from HSH1 and HSH2 of H3, aligned and displayed as in Fig. 2. The two helices in HSH1 are more nearly perpendicular to each other than those in HSH2.

hydrogen-bonding residues in the protein interior, but rather are subunit-assembling or DNA-binding residues.

A visual comparison (Fig. 2) of the structures of the four core histones with the globular portion of the linker histone H1 (GH1) (3, 8) demonstrates that their involvement in DNA binding and compaction is not the result of a shared structural ancestor. Overall, the size and shape of the two types of structure are quite different.

The Fold as a Module of Nucleosome Assembly. The histone fold is engaged directly in the formation of the histone dimers (1) and specifies the paired-element motifs (PEMs) that guide the docking of the DNA to the octamer (7). It is well established that from all possible pairwise associations of the histones certain histone dimers appear to be favored *in vivo* and/or *in vitro* (9, 10). The principles guiding the selection of "favored" partners may be responsible for the apparent limited sequence divergence at these, mostly hydrophobic, histone dimerization interfaces and for the overall evolutionary stability of histones.

However, during the life cycle of chromatin the nucleosome engages in a variety of structure/function transitions, most of which are expected to derive from changes in octamer structure. We believe that such modulations could be facilitated by structure diversification at the level of the histone dimers. Indeed, such limited diversification occurs at the areas of dimer-dimer contacts (tips of dimers) involved in generating the protein superhelix of the octamer. The two H3-H4 dimers associate and form the (H3-H4)₂ tetramer, and the contacts of the two opposing H3s (the HSH2 motifs of H3) are mainly hydrophobic in character (1). The analogous areas involved in the (H3-H4)/(H2A-H2B) associations are less hydrophobic and, consequently, this dimer-tetramer interface can be easily modulated by subtle environmental perturbations (11). The analogous portion of H2A is considerably less homologous to H3, and the tip of the dimer defined by this H2A segment

constitutes the end of the histone superhelix within the octamer and is most likely responsible for the "capping" (12) phenomenon in octamer assembly.

Dimer diversification also influences the way DNA interacts locally with the histones. The angle between helices I and II in H2A is different from that seen in the other core histones. This alteration in structure is a clear example of a single structural element linking two functional processes—i.e., octamer assembly and octamer-DNA binding. First, the carboxyl end of helix I and its adjacent loop form the interface between the two H2A/H2B dimers, an interface that may well contribute to the positive cooperativity observed during octamer assembly (11). Second, while the pitch of the protein superhelix becomes significantly steeper in the H2A-H2B domains, the change in pitch seen by the nucleosomal DNA as it passes over these areas is smoothed by the altered position of the amino end of helix I, which provides an intermediate docking pad between the (H3-H4)₂ tetramer and the H2A-H2B dimer.

The variations between the HSH1 and HSH2 segments of the fold appear to underscore differences in the way each histone dimer interacts with other dimers and with DNA. In general, the HSH2 motifs are more tightly conserved than the HSH1 motifs (Table 2), consistent with the fact that in the formation of the protein superhelix of the octamer, HSH2 enters into more protein-protein contacts than HSH1. Furthermore, in HSH2 the loop between helix II and strand B contains the well-conserved Lys-Arg pair that appears to be important in DNA binding. Two distinct regions from HSH1 in every fold—i.e., the amino terminus of helix I and strand A—interact with consecutive turns of the DNA and contact two different strands of the double helix separated by the width of the major groove. Therefore, the separation and relative orientation of helix I and strand A are fixed by the requirement that the HSH1 motif maintain the correct stereochemical correspondence to form a partnership with partly dehydrated DNA (13) curved at a radius appropriate to coil tightly around the octamer. Under the combined load of these structural requirements and their linkage to such indispensable functions as replication and transcription, the pressure for high stringency of conservation and simultaneously limited divergence of the fold characteristics during the evolution of the histone gene is easily appreciated.

Universal Occurrence of the Fold. Histones are universally found in animals, plants, and lower eukaryotes and have been recognized as the mediators of DNA compaction into chromatin (14). However, recently, small histone-like proteins (HMf, HMt) have been isolated from two strains of archaeobacteria and are reported to form dimers (ref. 15 and references therein) and to induce DNA supercoils (16). Two of these sequences (HMf_B and HMt_B) are shown in Fig. 1. These proteins are only 68 amino acids long, but, by homology, they appear to consist of a histone fold with two non-fold residues at both ends, and no labile amino termini. It is noteworthy that the HMf sequences possess a higher homology to the fold region of each core histone (24–29%) than the core histones display toward each other. Based on this high degree of homology we have modeled these archaeobacterial sequences within the eukaryotic histone fold (Fig. 4). The usually hydrophobic pattern required for pairwise interactions in eukaryotic histones is maintained in HMf and HMt, as is the pattern of positively charged residues required for DNA binding. On the basis of structural criteria they appear as, and we believe they are, true histones. However, because they are considerably smaller than the core histones and possess very short, if any, labile termini, the archaeal histones would therefore lack sites for post-translational modifications analogous to eukaryotic histones.

Table 2. Pairwise comparisons of histone fold amino and carboxyl halves

	H2A	H2B	H3	H4
H2A		0.9	0.7	1.0
H2B	1.8		1.0	0.7
H3	2.2	1.0		1.0
H4	2.1	1.7	1.7	

The distances shown represent the rms deviation (in Å) of α carbon positions between the two histones indicated. Below the diagonal are the comparisons for HSH1, calculated for residues 1–13 and 16–33 (residues 14 and 15 ignored because of the H4 deletion). Above the diagonal are the comparisons for HSH2, calculated for residues 34–65. Computations were performed as in Table 1.

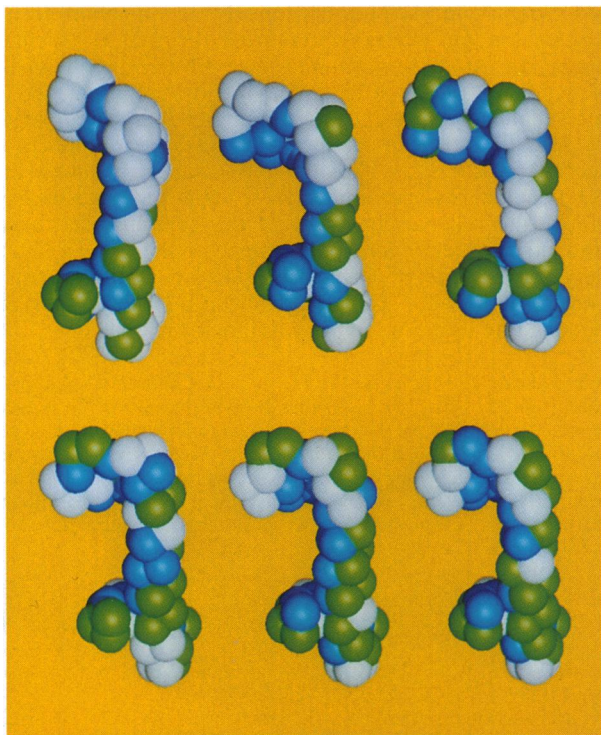


FIG. 4. Patterns of distribution of homologous residues within the fold. The chicken and the archaeal sequences (see Fig. 1) are presented in space-filling format. Amino acids have been lumped into three broad categories—polar (green), neutral (white), and strongly hydrophobic (blue)—based on their hydrophobicity as calculated by Eisenberg and McLachlan (17). The chicken histone structures are based on α carbon information from ref. 1. The HMFb and HMTb images have been derived by homology-based rendering relative to the H4 structure. From left to right: top row, H2A, H2B, H3; bottom row, H4, HMFb, HMTb. This display was generated by MIDAS (18), using twice the standard van der Waals radius for α carbons.

DISCUSSION

Evolutionary Aspects of the Fold. The ubiquitous utilization of the histones for the compaction of the genetic material suggests that from very early times, nature has selected this motif as the fundamental structural element for the reversible compaction of DNA. In search of better insights into the significance and potential function of this protein motif, we present here a close examination of the structural and evolutionary characteristics of the histone fold.

The two most salient features of the fold are (i) the twofold repetition of the helix/loop-and-strand/helix (HSH) configuration and (ii) the high degree of helicity (75%) that dominates the secondary structure of the fold (1). In Fig. 5 we illustrate the utilization of the histone fold elements in the formation of the paired element motifs (7) that serve as docking pads in DNA binding. Since the HSH motif is seen twice per histone and is present in all four core histone classes, it emerges as the basis from which eight classes of successful variations on the original motif evolved over time. It appears that evolution allowed considerable variation in primary structure, but only to the extent that the pattern of the histone fold was preserved. We now examine the possible significance of this conservation.

Within the nucleosome, many histone residues are involved in critical contacts with other histones that are necessary for the maintenance of the shape that is relevant to other chromosomal molecules. If, for the sake of argument, each of these contact domains is considered an "active site" indispensable to the function of the nucleosome, then it becomes obvious that the histone octamer is subject to multiple and simultaneous

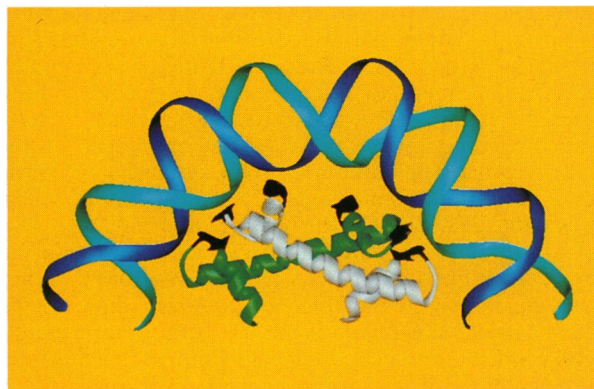


FIG. 5. A histone dimer is formed by the head-to-tail association of the fold portion of two chains, here, H3 (green) and H4 (white). The pseudo-twofold axis that relates H3 to H4 is in the plane of the paper (from top to bottom). The pairing of the two folds generates a smoothly curving outer surface containing the three PEMs (black) that dock to the inner face of nucleosomal DNA. This figure was generated by MIDAS (18).

selection pressures. Thus, there are few, if any, evolutionary "neutral" residues within the octamer—i.e., residues not involved in crucial histone–histone or histone–DNA contacts, or contacts with other regulatory elements. The histone fold appears to have been selected to fulfill most of these functions, especially the primary compaction of DNA. Additional and differential DNA compaction is derived from the extra-fold histone elements. We propose that the overall configuration of the fold within the octamer is strictly maintained through evolution by the requirement that three well-separated regions of the fold (docking pads) be spaced so as to interact with three consecutive turns of the phosphate backbones of a tightly curved double helix (7)—i.e., to bring nucleosome formation to a thermodynamic optimum. This requirement is met so precisely by the eight individual histone folds in the octamer that only a few (± 5) degrees of variation are seen over the 145° subtended by the superhelical surface of each histone dimer (Fig. 5). The relative positions and characteristic stereochemistry of the DNA docking pads make the architecture of the fold an excellent nonspecific DNA-binding protein motif.

The existence of the highly conserved histone fold offers strong evidence that the histones evolved from a common "protohistone" ancestor. The presence in archaeobacteria of histone-like proteins, each of which is more similar to any of the mammalian core histones than any one of the core histones is to the others, argues strongly for an early common ancestor. The packing of the eight chains in the histone octamer offers further evidence for a protohistone ancestor. When one histone assembles with its partner to form a dimer, their fold portions are related by a pseudo-twofold axis, a compelling argument for a single ancestor (6). The attributes of the histone fold and its utilization in dimer formation lead us to propose that the ancestral protohistone might have formed homodimers, with one chain being related to the other by a true twofold axis in that case. Such homodimers could assemble in a homotypic octamer with nucleosome-forming properties. However, such a homotypic octamer would have had quite different regulatory characteristics since, in the absence of accessory capping factors, it would also have had the potential to form long and inflexible homopolymers by end-to-end association between its symmetric dimer subunits.

Many of the residues in an individual histone are utilized in DNA binding, but in the fold region of each histone the double helix makes extensive contacts with the protein surface at the PEMs (7) (docking pads) of each dimer. Residues that can bind to B-form DNA in a sequence-independent fashion are found

in those areas—the two loops, the two β bridges, and the amino end of helix I—of the fold. In some pads, this binding always involves a positively charged residue, lysine or arginine, but other locations seem to tolerate a more general type of interaction, as implied by the different sizes and types of hydrogen-bonding residues present there, often serine, threonine, or tyrosine. The differences in the stereochemical properties of these residues may well reflect local and physiologically essential differences in the binding and bending of DNA. Finally, some portions of the fold are external or on the surface of the octamer but are not DNA-binding. We note that these regions are largely nonhomologous, although they are well conserved in individual types of histones, and perhaps underlie differing, albeit essential, but not yet identified functions such as interactions with transcription factors or participation in higher-order chromatin structures.

Specialized, Histone-Fold-Containing Proteins. Earlier, we utilized the amino acid sequences defined by the fold structures of the core histones to generate a consensus histone fold probe and used it to search a large set of protein sequence data (2). We found this fold sequence present in various proteins, several of which were not previously considered related to histones. Here we present a few examples from that search.

MacroH2A (19), a protein isolated from chromatin and nucleosomal structures, contains a full H2A sequence at its amino-terminal region and is contiguous with a leucine zipper sequence. No function has been identified yet for this protein. CENP-A is a centromere-specific protein identified as a component of nucleosomes and contains stretches of sequence highly similar to the fold region of H3, including nearly 70% identity with most of the HSH2 motif (20). The histone fold of CENP-A has been found to be required for targeting this protein to the centromere (21). Although many of the sequence alterations in CENP-A are conservative, two sets have strong functional implications. First, two inserted residues in CENP-A are likely to create extra bulk on the flat side of the octamer wedge, the side that might be involved in inter-octamer contacts. Second, the unstructured 42 residues of its amino terminus have almost no sequence similarity with the canonical H3 sequences. Both of these alterations occur in areas implicated in higher-order structure and thus are good candidates for generating changes in the centromeric domains of chromosomes.

Amino acid sequences characteristic of the histone fold have also been found in two proteins that are subunits of the *Drosophila* transcription initiation factor TFIID (2, 22). Although it is not known whether p42 and p62 form dimers, tetramers, or multimers, it is interesting to note that the HSH2 motif, which forms the contact surface between pairs of dimers in the histone octamer, is almost exactly duplicated in p42 and p62. Therefore, not only do these subunits have the potential to form a tetramer, but we predict that they also could form a heterotetramer with an H3–H4 pair.

In conclusion, the histone fold motif emerges as a well-preserved element of evolution of protein structure from bacteria to man. In the histones, it has been diversified to provide for the assembly of an oligomeric (octameric) articu-

lated protein endoskeleton (7) for DNA compaction. In this endoskeleton the central segment of the fold is the main element of histone dimerization. The differentiations of the end segments of the fold define mainly the properties of dimer–dimer contacts and the capping of the protein superhelix at the level of the octamer. However, the multiplicity of selective pressures for simultaneous octamer assembly and DNA docking may have limited these variations and thus have fostered the relative constancy of the core histone chains. Although first identified in core histones, the histone fold now emerges as a general protein-dimerization motif present in several proteins with specialized functions, such as enzymes, DNA-binding proteins, and transcription factors.

We dedicate this paper to the memory of Christian B. Anfinsen, whose inspired book “The Molecular Basis of Evolution” introduced one of us (E.N.M.) to the issues of macromolecular folding and evolution and has shaped the course of this research. We particularly appreciated his daily support and encouragement that sustained us during the last few years. This work was supported, in part, through a gift from Dr. G. Scangos of Bayer Corporation and by kind donations of several friends.

1. Arents, G., Burlingame, R. W., Wang, B.-C., Love, W. E. & Moudrianakis, E. N. (1991) *Proc. Natl. Acad. Sci. USA* **88**, 10148–10152.
2. Baxevanis, A. D., Arents, G., Moudrianakis, E. N. & Landsman, D. (1995) *Nucleic Acids Res.* **23**, 2685–2691.
3. Ramakrishnan, V., Finch, J. T., Graziano, V., Lee, P. L. & Sweet, R. M. (1993) *Nature (London)* **362**, 219–223.
4. Kraulis, P. J. (1991) *J. Appl. Crystallogr.* **24**, 946–950.
5. Chothia, C. & Lesk, A. M. (1986) *EMBO J.* **5**, 823–826.
6. Creighton, T. E. (1983) *Proteins; Structures and Molecular Properties* (Freeman, New York).
7. Arents, G. & Moudrianakis, E. N. (1993) *Proc. Natl. Acad. Sci. USA* **90**, 10489–10493.
8. Bernstein, F. C., Koetzle, T. F., Williams, G. J. B., Meyer, E. F., Jr., Brice, M. D., Rodgers, J. R., Kennard, O., Shimanouchi, T. & Tasumi, M. (1977) *J. Mol. Biol.* **112**, 535–542.
9. Pehrson, R. & Bustin, M. (1975) *Biochemistry* **14**, 3322–3331.
10. D’Anna, J. A. & Isenberg, I. (1974) *Biochemistry* **13**, 4992–4997.
11. Eickbush, T. H. & Moudrianakis, E. N. (1978) *Biochemistry* **17**, 4955–4964.
12. Baxevanis, A. D., Godfrey, J. E. & Moudrianakis, E. N. (1991) *Biochemistry* **30**, 8817–8823.
13. Eickbush, T. H. & Moudrianakis, E. N. (1978) *Cell* **13**, 295–306.
14. van Holde, K. E. (1988) *Chromatin* (Springer, New York).
15. Tabassum, R., Sandman, K. M. & Reeve, J. N. (1992) *J. Bacteriol.* **174**, 7890–7895.
16. Musgrave, D. R., Sandman, K. M. & Reeve, J. N. (1991) *Proc. Natl. Acad. Sci. USA* **88**, 10397–10401.
17. Eisenberg, D. & McLachlan, A. D. (1986) *Nature (London)* **319**, 199–203.
18. Ferrin, T. E., Huang, C. C., Jarvis, L. E. & Langridge, R. (1988) *J. Mol. Graphics* **6**, 13–27.
19. Pehrson, J. R. & Fried, V. A. (1992) *Science* **257**, 398–400.
20. Palmer, D. K., O’Day, K., Trong, H. L., Charbonneau, H. & Margolis, R. L. (1991) *Proc. Natl. Acad. Sci. USA* **88**, 3734–3738.
21. Sullivan, K. F., Hechenberger, M. & Masri, K. (1994) *J. Cell Biol.* **127**, 581–592.
22. Kokubo, T., Gong, D.-W., Wootton, J. C., Horikoshi, M., Roeder, R. G. & Nakatani, Y. (1994) *Nature (London)* **367**, 484–487.