

Modern proteomes contain putative imprints of ancient shifts in trace metal geochemistry

Christopher L. Dupont, Song Yang, Brian Palenik, and Philip E. Bourne

PNAS published online Nov 10, 2006;

doi:10.1073/pnas.0605798103

This information is current as of February 2007.

Supplementary Material

Supplementary material can be found at:
www.pnas.org/cgi/content/full/0605798103/DC1

This article has been cited by other articles:
www.pnas.org#otherarticles

E-mail Alerts

Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).

Rights & Permissions

To reproduce this article in part (figures, tables) or in entirety, see:
www.pnas.org/misc/rightperm.shtml

Reprints

To order reprints, see:
www.pnas.org/misc/reprints.shtml

Notes:

Modern proteomes contain putative imprints of ancient shifts in trace metal geochemistry

Christopher L. Dupont^{*†}, Song Yang[‡], Brian Palenik[§], and Philip E. Bourne^{¶||}

Departments of [†]Chemistry and Biochemistry and [¶]Pharmacology, ^{||}San Diego Supercomputer Center, and [§]Marine Biology Research Division, ^{*}The Scripps Institution of Oceanography, University of California at San Diego, La Jolla, CA 92093

Edited by Janet M. Thornton, European Bioinformatics Institute, Cambridge, United Kingdom, and approved September 25, 2006 (received for review July 10, 2006)

Because of the rise in atmospheric oxygen 2.3 billion years ago (Gya) and the subsequent changes in oceanic redox state over the last 2.3–1 Gya, trace metal bioavailability in marine environments has changed dramatically. Although theorized to have influenced the biological usage of metals leaving discernable genomic signals, a thorough and quantitative test of this hypothesis has been lacking. Using structural bioinformatics and whole-genome sequences, the Fe-, Zn-, Mn-, and Co-binding metallomes of 23 Archaea, 233 Bacteria, and 57 Eukarya were constructed. These metallomes reveal that the overall abundances of these metal-binding structures scale to proteome size as power laws with a unique set of slopes for each Superkingdom of Life. The differences in the power describing the abundances of Fe-, Mn-, Zn-, and Co-binding proteins in the proteomes of Prokaryotes and Eukaryotes are similar to the theorized changes in the abundances of these metals after the oxygenation of oceanic deep waters. This phenomenon suggests that Prokarya and Eukarya evolved in anoxic and oxic environments, respectively, a hypothesis further supported by structures and functions of Fe-binding proteins in each Superkingdom. Also observed is a proliferation in the diversity of Zn-binding protein structures involved in protein–DNA and protein–protein interactions within Eukarya, an event unlikely to occur in either an anoxic or euxinic environment where Zn concentrations would be vanishingly low. We hypothesize that these conserved trends are proteomic imprints of changes in trace metal bioavailability in the ancient ocean that highlight a major evolutionary shift in biological trace metal usage.

bioinorganic chemistry | evolution | fold families | structural bioinformatics

The emergence of oxygenic photosynthesis is associated with major changes in global biogeochemistry and metabolism (1, 2). In particular, the rise in atmospheric oxygen \approx 2.3 billion years ago (Gya) (3, 4) potentially led to the oxygenation of the entire ocean (5), whereas an alternative theory proposes that the deep ocean became euxinic (anoxic and sulfidic) \approx 1.8 Gya (6, 7), before an oxygenation of deep waters \approx 1 Gya (8). Putting aside for now when and where, these changes in the overall redox state of the ocean would dramatically influence trace metal chemistry and bioavailability, with an anoxic ocean being characterized by relatively high Fe, Mn, and Co but low Zn concentrations (9) (Fig. 4, which is published as supporting information on the PNAS web site). A euxinic ocean would have comparatively lower concentrations of all of these metals, particularly Zn (9) (Fig. 4). The oxygenation of oceanic deep waters would have dramatically increased Zn concentrations, with concomitant yet less severe decreases in Fe, Mn, and Co levels (9) (Fig. 4). As postulated by Williams and Frausto da Silva (10), these drastic shifts in metal bioavailability theoretically influenced the selection of trace elements for biological usage, leaving a record within the genomes and proteomes of extant organisms.

Protein structure has a remarkable level of redundancy, with a limited number of 3D folds describing all of life (11). Further, structure is retained over long evolutionary time scales, even

when most sequence homology is lost, providing an excellent tool for this study. Already, the identification of domains within protein structures and the systematic and hierarchical classification of these domains have been used to study evolution (12). Within these hierarchical classifications reside fold superfamilies (FSF) and fold families (FF); a FSF contains structures believed to be evolutionarily related despite a lack of clear sequence similarity, whereas a FF contains structures with evident structural, functional, and sequence similarities (a FSF is composed of one or more FF). The gain or loss of a FSF or FF by an organism constitutes an important evolutionary event, either reducing or expanding the repertoire of functions available to that organism. Indeed, the presence or absence of FSFs in a proteome has been shown to discriminate species well enough to construct reasonable phylogenetic trees for all of life (13).

Here, we used structural bioinformatics to study the distribution of metal-binding protein structures within the proteomes of Archaea, Bacteria, and Eukarya, which to our knowledge has not been done before. The results suggest that ancient changes in trace metal geochemistry do indeed leave imprints observable within the genomes and proteomes of modern life and provide an important constraint on the evolution of Eukarya.

Results and Discussion

The Superfamily database (14, 15), derived from the Structural Classification of Proteins (SCOP) (16), provides an independent assessment of the presence and abundance of structural domains belonging to FSFs and FFs across a diverse set of species for which complete genome and translated proteome sequences are available. To extract the desired information from the Superfamily database, we manually annotated SCOP version 1.69 according to metal binding (Tables 3–6, which are published as supporting information on the PNAS web site). Both the raw structural data from the Protein Data Bank (PDB) (17) and the primary literature associated with each structure were used to identify covalently bound metals. Protein domains that bind a metal-containing cofactor (e.g., Co-containing B₁₂ and Fe-containing heme) were considered metal binding. Here, an ambiguous FSF is defined as one in which the structures comprising that FSF bind different metals or contain a combination of both metal- and nonmetal-binding structures. Likewise, an ambiguous FF contains a mixture of metal- and nonmetal-binding structures or structures binding different metals. Approximately half of the metal-binding FSFs and 10% of the metal-binding FFs are ambiguous (Table 7, which is pub-

Author contributions: C.L.D., B.P., and P.E.B. designed research; C.L.D. and S.Y. performed research; C.L.D. analyzed data; and C.L.D., B.P., and P.E.B. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS direct submission.

Freely available online through the PNAS open access option.

Abbreviations: FSF, fold superfamily; FF, fold family; SCOP, structural classification of proteins; PDB, Protein Data Bank; Gya, billion years ago.

[†]To whom correspondence should be addressed. E-mail: cdupont@ucsd.edu.

© 2006 by The National Academy of Sciences of the USA

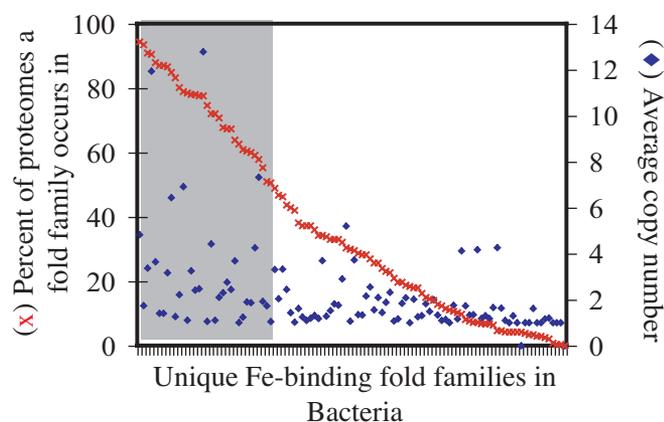


Fig. 2. Diversity and abundance of Fe-binding fold families in Bacteria. For each Fe-binding fold family (tick marks on x axis), the red \times (left axis for scale) shows the percentage of proteomes in which it occurs, whereas the blue \blacklozenge (right axis for scale) shows the average copy number in proteomes where it does occur. The shaded area highlights the number of fold families that occur in at least 50% of the Bacterial proteomes examined. Similar trends are observed for the other metals and Superkingdoms.

web site), but more complete proteomes are needed to robustly test this relationship.

The observed scaling is not due to a core set of abundant metal-binding proteins (e.g., a metal-binding superfold) within a given Superkingdom; individual species have drawn broadly and diversely from the pool of available metalloproteins. That is, very few metal-binding domains are ubiquitous or are found within all of the proteomes of a Superkingdom, yet many are present in at least one proteome (e.g., see Fig. 2 for the case of Fe in Bacteria). Additionally, structural domains found in all or most of the proteomes are not necessarily more abundant in those proteomes than structural domains found in only a few proteomes (Fig. 2). Essentially, different organisms have different metal-binding domains, a logical extension of the results of Yang *et al.* (13), yet the total abundances of Fe-, Mn-, Zn-, or Co-binding domains within a proteome conform to fundamental constants defined by power laws.

It appears that methodological limitations, including proteome coverage and sampling bias, do not contribute to the observed trends. According to Superfamily (14), on average 55% of Archaeal and Bacterial proteomes and 40% of Eukaryotic proteomes have fold families assigned. Although this coverage may seem limiting, results from the Protein Structure Initiative suggest that $\approx 90\%$ of this unannotated space is comprised of variants of already discovered fold families (20), and that only 10% of the undiscovered fold families will actually be metal-binding (21). It appears that membrane proteins are similarly distributed, with the most abundant membrane folds being already described (22). Essentially, it seems unlikely that a new protein fold family will be discovered that is abundant enough to overly skew the observed results.

A further concern is that the observed results are biased by the available whole genomes. The Archaea sequenced are mostly thermophiles from anoxic environments (although this potentially provides a modern-day glimpse into ancient bioinorganic chemistry), whereas the sequenced Eukarya are almost entirely aerobic. The dataset does include the Eukaryotic anaerobic amitochondriotic parasite *Encephalitozoon cuniculi*, which has metallomic features typical of aerobic Eukaryotes. In contrast, the analyzed Bacteria are from a broad array of environments and have a variety of oxygen tolerances and therefore can be used to gain a preliminary understanding of the influence of modern environment and metabolism on metallomic content.

Surprisingly, within the Bacterial Superkingdom, differences in oxygen tolerance do not seem to influence the proteomic abundance of metal-binding domains (Fig. 8, which is published as supporting information on the PNAS web site). Instead, Phylum or Classes roughly group together (Fig. 9, which is published as supporting information on the PNAS web site), implying that the observed stoichiometries are vertically inherited. Future work will explore the causes of the nonsize-dependent variance within the data ($<10\%$; see Table 1). Additionally, more genome and proteome sequences will allow for a continued updating of these results, with the proteomes of aerobic Archaea and anaerobic Eukaryotes providing key tests.

Accepting that the data are not limiting, the critical question remains as to the source of the Superkingdom level differences in the power law slopes. It seems reasonable to assume that the observed power law slopes are determined by selective pressure (23), and that trace metal bioavailability can produce such pressure (24). Hence, we hypothesize that the environmental bioavailability of trace elements during major periods of phylogenetic diversification shape the evolution of vertically inherited metal homeostasis systems that then continually influence the retention and loss of genetic material. The differences in the power law scaling for metal-binding structures within Prokaryotic and Eukaryotic proteomes are similar to the shifts in trace metal bioavailability caused by increasing oxygen, implying that Prokaryotic and Eukaryotic organisms diversified in anoxic and oxic environments, respectively. The proposed theory entails a closed feedback loop, whereby a biological phenomenon (cyanobacterial production of oxygen) incites a shift in trace metal geochemistry that in turn influences the evolution of bioinorganic chemistry. An alternative theory is that the observed differences between the Prokaryotes and Eukaryotes are due to an unknown but environmentally unrelated phenomenon. To address this possibility, we further examined the functions and structures of Zn- and Fe-binding proteins for environmentally consistent signals.

As stated, Eukaryotic and Prokaryotic proteomes show significant differences in the abundance of Zn-binding domains. These are wholly attributable to structures in the “small protein” structural class (Fig. 3A), typified by small Zn-binding domains such as Zn fingers and RING domains involved in protein–DNA/RNA interactions and protein–protein interactions, respectively. Eukaryotic proteomes also encode for a greater structural diversity of “small protein class” fold families that bind Zn (Fig. 3B), a noteworthy radiation in the diversity and usage of Zn within proteins, one that is predominantly structural. Most of the “small protein” class Zn-binding protein fold families are unique to Eukarya, although a significant subset is shared with Archaea (Fig. 3B). Because Zn concentrations would be vanishingly low in an anoxic or euxinic environment (ref. 9; Fig. 4), it seems unlikely that such a diversification in the biological usage of Zn could occur under such conditions.

The functions and structures of the prevalent Fe-binding domains in each Superkingdom also are consistent with the evolution of Eukarya in an oxic environment. Fe-binding FFs were characterized according to the mode of Fe binding (Fe-S, heme, or direct amino acid), and the abundances of these binding forms were quantified for each Superkingdom (Table 2). Archaeal and Bacterial metallomes have significantly more Fe-S proteins and fewer heme proteins than the Eukaryotic metallomes (Table 2). Both the observed Fe-S clusters and hemes function in e^- transfer reactions, but Fe-S clusters are oxygen-sensitive and have more negative reduction potentials than heme-based Fe proteins such as cytochromes (25). The proteomes of aerobic Bacteria also contain fewer Fe-S clusters and more hemes than anaerobic Bacteria (Table 9, which is published as supporting information on the PNAS web site), suggesting that the actual repertoire of metalloproteins within the con-

Eukarya to the late Proterozoic (0.9–1.2 Gya; ref. 29). The deep ocean was potentially anoxic or euxinic during both of these periods (30); these data have been used by some to argue that Eukarya evolved and diversified in anaerobic environments (31). This contention is contrary to the theory proposed by Anbar and Knoll that low Cu and Mo bioavailability in a euxinic ocean limited Eukaryotic diversification (30). Our results support the latter hypothesis that oxygen-induced changes in trace metal bioavailability occurred before the diversification of Eukarya and implicate Zn as another relevant metal. The fossil record indicates early evolutionary radiations of Eukarya likely occurred in shallow coastal environments (32, 33), where a combination of high oxygen concentrations and a terrestrial supply of trace metals may have increased the bioavailability of Zn, Cu, and Mo. Note that, although oxic microenvironments may have existed in the surface ocean since the advent of oxygenic photosynthesis, the supply of trace metals to these microenvironments would have been unchanged until the oxygenation of deep waters, in contrast to coastal environments.

The idea that the rise in oxygen affected the usage of trace metals was originally proposed by Williams and Frausto da Silva (10), and a few studies have used sequence-based methods to study the coevolution of biology and geochemistry. Morgan *et al.* (34) found that Eukarya have a higher diversity and abundance of Ca^{+2} -binding protein sequence families. Zerkle *et al.* (35) examined the distributions of ORFs annotated as known metal-binding proteins within the genomes of Prokarya, finding differences based upon metabolism and phylogeny. The analysis conducted here expands on these theories and efforts. Within a proteome, the abundances of metal-binding domains conform to a stoichiometry defined by evolutionary constants despite the wide diversity of physiologies and environments of the analyzed organisms. Further, these constants exhibit Superkingdom-specific behavior consistent with development within anoxic vs. oxic environments. It must be noted that the observed proteomic stoichiometries likely do not define the physiological metal requirements of a specific organism. Single metalloproteins can constitute a large portion of an organism's metal usage. However, whole proteomes are less susceptible to gene acquisition events or evolutionarily recent ecological or physiological adaptations, such as those observed in coastal and open ocean cyanobacteria (36). Hence, we feel that the whole-proteome patterns observed represent a broader and more durable view into the ancient environment of Earth than physiological quotas.

Methods

Data Sources. SCOP (16) provides a hierarchical classification of all protein domains published in the PDB (17). SCOP Version 1.69 has sorted 70,800 domains into 945 defined folds that are assigned to 1,539 superfamilies and further subdivided into 2,845 families. The Superfamily database (14, 15) was the source of all domain assignments; Release 1.69 covers 313 complete genomes (23 Archaea, 233 Bacteria, and 57 Eukaryota). The Superfamily database, using a hybrid approach of a hidden Markov model searching protocol and subsequent pairwise comparisons (15), uses a probability cutoff of $E = 2 \times 10^{-2}$ for identifying likely members of a group; it also provides a confidence level (in the form of an E value) for every candidate identified. As was done by Yang *et al.* (13), a more stringent E value cutoff of 10^{-4} was used for the domain assignments here.

Annotation of SCOP per Metal Binding. Each FSF in the SCOP database was manually examined for structures containing a covalently bound inorganic ion. This objective requires examining the FF, fold domains, and specific example structures within each FSF. For a FSF or FF to be considered metal-binding, only one of the representative structures has to contain a bound metal. If all of the representative structures in a FSF or FF bind the same metal, the family is considered unambiguous; only these FFs were used for this study.

Automated annotations of SCOP in this fashion are inhibited by two factors: (i) Some structures are crystallized with nonnative metals, and (ii) some PDB data files are less thorough in the description of binding mode and domain. The manual examination procedures were simplistic; the accompanying PDB file was examined for an inorganic ion and covalent binding of that ion. Some PDB files provide metal-binding information (i.e., heme, amino acid, and specific binding residues), and whether the metal is native or simply part of the crystallization buffer. In the cases where this was not clear, the primary literature citation was examined. Attention was also paid to which structural domains actually bind the metal. For example, there are numerous distinct FSFs and FFs that contain domains of cytochrome *c* oxidase in the SCOP, and each entry states “complexed with cdl, chd, cu, cua, dm, hea, mg, na, pek, pgv, psc, tgl, unx, zn,” yet only a select few are Cu-, Fe-, Mg-, or Zn-binding. The FFs that unambiguously bind Fe, Zn, Mn, or Co are shown in Tables 3–6.

Data Management and Analysis. Matrices were constructed, with each row representing a distinct species and each column representing the abundances of a specific metal-binding FF in each species. For the power law distributions in Fig. 1, the FFs in a given proteome that unambiguously bind a metal were summed and plotted against the total number of domains assigned to that proteome (the sum of all FF assignments for a proteome). For Table 2, the matrices were normalized by dividing the abundances of each FF within a species by the total number of structural domains assigned to that species' proteome. These internal percentages were then averaged over the entire Superkingdom.

Power law fits were determined in Matlab (Mathworks, Natick, MA). The data were log-transformed, and the linear fit was found by using a geometric mean least-squares fitting technique. Groupings of points were compared by using ANCOVA (the *aocool* function in Matlab). Slopes were compared by using the *multicompare* function. Power law fit qualities (F values; Table 1) were determined by using the method of van Nimwegen (18). Briefly, the data were log-transformed. The distance from each point (defined by x_i, y_i) to the center of the scatter (D_c) is determined by $D_c = \sqrt{((x_i - \text{mean of all } x)^2 + (y_i - \text{mean of all } y)^2)}$. Then the distance of each point to the fitted line (D_l) is determined by $D_l = \sqrt{((y_i - mx_i - b)^2 / (m^2 + 1))}$, where m and b equal the power law slope and intercept. Then the fraction of the variance explained by the data is given by $F = 1 - (\sum(D_l)^2) / (\sum(D_c)^2)$.

We thank K. N. Chang, A. J. Lucas, R. J. P. Williams, S. Veretnik, and anonymous reviewers for constructive comments and suggestions. P.E.B. was supported by National Institutes of Health Grant GM63208, and C.L.D. is grateful for funding from a National Defense Science and Engineering Graduate Research Fellowship and the National Science Foundation/Department of Energy-supported Princeton Center for Bioinorganic Chemistry.

1. Kopp RE, Kirschvink JL, Hilburn IA, Nash CZ (2005) *Proc Natl Acad Sci USA* 102:11131–11136.
2. Raymond J, Segre D (2006) *Science* 311:1764–1767.
3. Bekker A, Holland HD, Wang PL, Rumble D, Stein HJ, Hannah JL, Coetzee LL, Beukes NJ (2004) *Nature* 427:117–120.
4. Farquhar J, Bao H, Thiemens M (2000) *Science* 289:756–758.

5. Holland HD (1984) *The Chemical Evolution of the Atmosphere and Oceans* (Princeton Univ Press, Princeton).
6. Canfield DE, Teske A (1996) *Nature* 382:127–132.
7. Arnold GL, Anbar AD, Barling J, Lyons TW (2004) *Science* 304:87–90.
8. Canfield DE (1998) *Nature* 396:450–453.
9. Saito MA, Sigman DM, Morel FMM (2003) *Inorg Chim Acta* 356:308–318.

