AlphaFold2-guided description of CoBaHMA, a novel family of bacterial domains within the Heavy-Metal Associated superfamily

Isabelle Callebaut¹, Geoffroy Gaschignard¹, Maxime Millet¹, Apolline Bruley¹, Karim Benzerara¹, Manuela Dezi¹, Feriel Skouri-Panet¹, and Elodie Duprat¹

¹Sorbonne Universite UFR de Physique

September 29, 2023

Abstract

Three-dimensional structure information, now available at the proteome scale, may facilitate the detection of remote evolutionary relationships in protein superfamilies. Here, we illustrate this with the identification of a novel family of protein domains related to the ferredoxin-like superfold, by combining (i) transitive sequence similarity searches, (ii) clustering approaches and (iii) the use of AlphaFold2 3D structure models. Domains of this family called CoBaHMA, were initially identified in relation with the intracellular biomineralization of calcium carbonates by Cyanobacteria. They are part of the large heavy-metal-associated (HMA) superfamily, departing from the latter by specific sequence and structural features. In particular, most CoBaHMA domains share conserved basic amino acids, forming a positively charged surface, which is likely to interact with anionic partners. CoBaHMA domains are found in diverse modular organizations in bacteria, existing in the form of monodomain proteins or as part of larger proteins, some of which are membrane proteins involved in transport or lipid metabolism. This suggests that the CoBaHMA domains may exert a regulatory function, involving interactions with anionic lipids. This hypothesis might have a particular resonance in the context of the compartmentalization observed for cyanobacterial intracellular calcium carbonates.

AlphaFold2-guided description of CoBaHMA, a novel family of bacterial domains within the Heavy-Metal Associated superfamily

Geoffroy Gaschignard^a, Maxime Millet^a, Apolline Bruley, Karim Benzerara, Manuela Dezi, Feriel Skouri-Panet, Elodie Duprat^b, Isabelle Callebaut^b

Sorbonne Université, Muséum National d'Histoire Naturelle, UMR CNRS 7590, Institut de Minéralogie, de Physique des Matériaux et de Cosmochimie, IMPMC, 75005 Paris, France

Geoffroy Gaschignard and Maxime Millet should be considered joint first authors.

Elodie Duprat and Isabelle Callebaut should be considered joint senior authors

Correspondence :

Isabelle Callebaut and Manuela Dezi, Sorbonne Universite, Museum National d'Histoire Naturelle, UMR CNRS 7590, Institut de Mineralogie, de Physique des Materiaux et de Cosmochimie, IMPMC, 75005 Paris, France. Email: isabelle.callebaut@sorbonne-universite.fr (I.C.) and Manuela.dezi@sorbonne-universite.fr (M.D)

Data availability statement

Supplementary Figures are included in Supporting Information (pdf file). Other data generated and used in this study are available as Supplementary Data and Table on Zenodo (DOI: 10.5281/zenodo.8387655). Scripts used in the present work are available in a GitHub repository at *https://github.com/GGasch/CoBaHMA_Detection*.

Funding Information

This work was supported by the Agence Nationale de la Recherche (ANR Harley, ANR-19-CE44-0017-01; ANR PHOSTORE, ANR- 19-CE01-0005).

Conflict of Interest statement: none declared.

Running title: The CoBaHMA family of domains

Acknowlegments

We warmly thank Jean-Paul Mornon for critical reading of the manuscript.

Supporting information

Supporting information can be found online in the Supporting Information section at the end of this article and on Zenodo (DOI: 10.5281/zenodo.8387655).

ABSTRACT

Three-dimensional structure information, now available at the proteome scale, may facilitate the detection of remote evolutionary relationships in protein superfamilies. Here, we illustrate this with the identification of a novel family of protein domains related to the ferredoxin-like superfold, by combining (i) transitive sequence similarity searches, (ii) clustering approaches and (iii) the use of AlphaFold2 3D structure models. Domains of this family called CoBaHMA, were initially identified in relation with the intracellular biomineralization of calcium carbonates by Cyanobacteria. They are part of the large heavy-metal-associated (HMA) superfamily, departing from the latter by specific sequence and structural features. In particular, most CoBaHMA domains share conserved basic amino acids, forming a positively charged surface, which is likely to interact with anionic partners. CoBaHMA domains are found in diverse modular organizations in bacteria, existing in the form of monodomain proteins or as part of larger proteins, some of which are membrane proteins involved in transport or lipid metabolism. This suggests that the CoBaHMA domains may exert a regulatory function, involving interactions with anionic lipids. This hypothesis might have a particular resonance in the context of the compartmentalization observed for cyanobacterial intracellular calcium carbonates.

KEYWORDS: Heavy-metal-associated, AlphaFold2, sequence similarity search, graph clustering, modular organization, functional annotation, P1B-ATPase, ABC transporter, PAP2, biomineralization, CoBaHMA

ABBREVIATIONS

aa : amino acids, A : Actuator, ABC:ATP-Binding Cassette, AF2 : AlphaFold2, AFDB : AlphaFold Protein Structure DataBase, ATPDB: ATP-Binding Domain,ccya: calcyanin, CoBaHMA: Conserved Basic residues HMA, DAG: diacylglycerol, GlyZip: glycine-zipper,HCA : Hydrophobic Cluster Analysis, HMA: Heavy Metal-Associated, MSA : Multiple Sequence Alignment, N:Nucleotide, NBD: Nucleotide-Binding Domain, P:phosphorylation, PAE: Predicted Aligned Error, PAP2:type 2 phosphatidylglycerol phosphatase, PDB: Protein Data Bank, PG : phosphatidylglycerol, pLDDT: predicted Local Distance Difference Test, R: regulatory, RSSs: Regular Secondary Structures, S: auxiliary membrane, SCOP e: Structural Classification of Proteins — extended, TM:Transmembrane, TMD : Transmembrane Domain

Introduction

Superfolds are folds observed in a large number of evolutionary unrelated protein domain superfamilies ¹ and are characterized by compact super-secondary structure patterns^{2,3}. One of these superfolds is the ferredoxin-like fold, found in 62 superfamilies according to the SCOPe classification (fold d.58, ⁴) and present in many domains with various functions ^{5,6}. It is made of a repeated β - α - β super-secondary structure, forming a four-stranded β -sheet, with the two α -helices packed into an α - β sandwich. As for other superfolds, the ferredoxin-like fold is subject to circular permutations, a mechanism which allows adaptation and the emergence of new functions. This can be visualized by the ligation of the amino- and carboxyl-termini and subsequent cleavage at another site ^{7,8}. Among the 62 superfamilies comprising the ferredoxin-like

fold, the heavy-metal-associated (HMA) superfamily (SCOPe d.58.17) contains the eponymous HMA family (**Figure 1-A**). Domains of this HMA family contain two conserved cysteine residues involved in heavy metal binding. They are found in a variety of metal-trafficking proteins, which play essential roles in transport and homeostasis 9,10 . While they can be found alone, HMA domains are also part of diverse domain architectures, especially associated with P1B-ATPases, which are integral membrane proteins allowing transport of metals across cell membranes 11 .

We recently identified a novel family of domains belonging to the HMA superfamily, called CoBaHMA (after Co nserved Ba sic residues HMA). This discovery originates from the characterization of a novel family of two-domain proteins, named calcyanin, which is associated with intracellular biomineralization of calcium carbonates by cyanobacteria¹². Calcyanins share a common architecture consisting in a conserved C-terminal domain, made of a three-fold repeated, unusually long glycine zipper motif $((GlyZip)_3)$. Glycine zippers themselves consist of repeated GXXXG motifs, commonly found in transmembrane domains and bacterial toxins 13,14 . The calcyanin (GlyZip)₃ domain is preceded by a variable N-terminal domain, specific to the 4 distinct calculation subgroups identified to this date, which are found in distinct clades of cvanobacteria. The N-terminal domain of one of these calcyanin subgroups is a CoBaHMA domain. It is represented by 15 sequences ¹². It shares significant sequence similarities with HMA domains, as well as plant integrated HMA domains ¹⁵ and YAM domains found in the C-terminal part of the bacterial YajR, an integral membrane protein which belongs to the major facilitator superfamily (MFS, ^{16,17}). These three families of domains share the same 3D structure, but only the first one (HMA) is referenced in the SCOPe classification (d.58.17.1, ⁴). Sequence analysis and molecular modeling ¹² have revealed specific features of the CoBaHMA domain family, including the presence of an additional strand $\beta 0$ at the N-terminus of the domain, which replaces the C-terminal β 4 strand, thus conforming to the circular permutation scenario described before (Figure 1-B). CoBaHMA domains in calcyanins are also characterized by the presence of conserved basic amino acids (aa) in $\beta 0$ (H) and $\beta 1$ (R/K) strands, forming a charged patch on one side of the β -sheet. Despite the clear identification of all these specific sequence and structural features, the function(s) of CoBaHMA domains remain(s) unknown.

Hosted file

image1.emf available at https://authorea.com/users/669589/articles/669968-alphafold2-guideddescription-of-cobahma-a-novel-family-of-bacterial-domains-within-the-heavy-metalassociated-superfamily

Figure 1: Topology diagrams and ribbon representations of the 3D structures of (A) an HMA domain (Human CopZ; experimental 3D structure - PDB 2QIF); (B) a CoBaHMA domain (S. calcipolaris calcyanin; AlphaFold2 model). Regular secondary structures are colored rainbow, from the N-terminus (blue) to the C-terminus (red), with the exception of the additional strand $\beta 0$, specific of the CoBaHMA domain (pink). Conserved amino acids, specific to the two families, are highlighted in a ball-and-stick representation.

Here, we searched for the presence of CoBaHMA domains in proteins distinct from calcyanins. We posit that the identity of co-occurring domains or combinations of domains in other proteins may inform about their functional context ¹⁸⁻²⁰. For this purpose, we propose a novel methodological framework, taking advantage of the structural information provided by the Alphafold2 (AF2) 3D structure models ²¹ of the retrieved sequences, now widely available in the Alphafold DataBase (AFDB) ²². This allows a both sensitive and specific detection of CoBaHMA domains within the HMA superfamily. As a result, we describe a wide diversity of modular organizations in which CoBaHMA domains are included, with some yet-to-be-characterized regions. Moreover, the CoBaHMA domain family appears to be specific to bacteria and frequently associated with transmembrane and soluble domains involved in transport of substrates and lipid metabolism, suggesting a regulatory role with regard to these functions, possibly via an interaction with charged lipids.

MATERIAL AND METHODS

The workflow developed and used in this study is shown in Figure 2.

Hosted file

image2.emf available at https://authorea.com/users/669589/articles/669968-alphafold2-guideddescription-of-cobahma-a-novel-family-of-bacterial-domains-within-the-heavy-metalassociated-superfamily

Figure 2: Workflow of the current analysis. (A) The 15 sequences of CoBaHMA domains were used as probes for sequence similarity searches with HHBlits against the Uniclust30 database. The 3D structure models of the sequences that were identified by HHBlits were considered in order to assess the presence of strand $\beta 0$, a feature allowing to distinguish true CoBaHMA domains. The new CoBaHMA domain sequences were then used as additional probes with HHBlits to perform transitive searches. (B) Full length proteins with at least one CoBaHMA domain were annotated relative to transmembrane regions, calcyanin signature, structural domains and domains families using deepTMHMM, PCALF, DomainMapper and InterProScan, respectively. In addition, the taxonomy of all sequences is retrieved through the UniProt API. (C) Sequences were clustered using mmseqs2 and representative sequences were searched against themselves. Resulting alignment scores were used to build a similarity network. The network was refined iteratively to ensure that all sequence within a community share a similar length with a maximal amplitude of 50 residues. Finally, nodes and edges were rendered using the edge-rendering and weighted spring embedded layout from Cytoscape.

Structure-guided detection of CoBaHMA domain sequences (Figure 2-A)

We searched the UniClust30/UniRef30 databases (version $2022_{-}02.$ downloaded from https://gwdu111.gwdg.de/~compbiol/uniclust/2022_02/)²³ using the standalone version of HHBlits²⁴. The 15 sequences of CoBaHMA domains identified by Benzerara et al. ¹² were used as individual probes. The sequence similarity search consisted in 8 iterations with hits covering at least 50% of the sequence length and an E-value threshold of 1e-3, for each sequence probe. For each iteration and each probe, we gathered the sequences extracted from the output a3m multiple sequence alignment. We removed trivial redundancy using an in-house python3 script, by merging entries with identical identifiers and whose sequences were identical on at least 50% of their respective length. In order to find which entry had a 3D model available in the AlphaFold database (AFDB) (https://alphafold.ebi.ac.uk/)²², we downloaded the accession_ids.csv file from http://ftp.ebi.ac.uk/pub/databases/alphafold/. Based on this file, we downloaded the 3D models corresponding to our entries from AFDB v4. We cropped each 3D model based on the boundaries of the CoBaHMA domain candidates identified with the HHBlits search. From the 3D coordinates of these sub-3D models, we computed the secondary structures associated with each aa using DSSP^{25,26}. The 8 states output of DSSP was converted to 3 states using the EVA convention²⁷. Following this convention, α -helix, 3_{10} helix and π -helix are converted to helices (H), extended and isolated β -bridge to extended (E), and turn, bend and other to coils (C). To validate the presence of a secondary structure, we set a threshold of at least 3 consecutive as with the same secondary structure assignment. As a result, we could infer the string of regular secondary structures of the sub-models. We set 4 criteria to identify a CoBaHMA domain, as follows. First, it should contain the secondary structure pattern EEHEEH(E) typical of CoBaHMA domains¹². Second, it should contain neither the Cys-X-X-Cys pattern, typical of HMA domains ¹⁰, nor the Cys-X-X-Cys pattern typical of TRASH domains ²⁸. Third, the secondary structure pattern should not start with EH, as observed in canonical HMA domains. Fourth, the first two beta strands should be separated by a coil of 1 or 2 aa, to avoid confusing a $\beta 0$ strand of a CoBaHMA domain with the $\beta 4$ strand of an HMA domain, which is followed by a second HMA domain starting with a $\beta 1$ strand. A final manual check based on the alignment of sequences was conducted to remove the remaining false positives. Overall, entries meeting these 4 criteria and passing the manual check were considered as CoBaHMA domains.

We performed transitive searches with newly identified CoBaHMA domains. This process was iterated 6 times. We removed the redundancies out of the 6 iteration outputs by merging the entries that had the same identifiers and shared at least 50% of their sequences.

Finally, since HHblits alignment only gave the portions of the sequences that match the probes, we down-loaded from UniProt (https://www.uniprot.org/)²⁹ the full sequences of each identifier that had at least one CoBaHMA domain in its sequence.

Annotation of the full-length sequences (Figure 2-B)

The TaxID and organism name associated with the sequences were retrieved using the UniProt REST API by downloading UniProt ttl files and serializing sequences information. Taxonkit ³⁰was used to convert TaxID into taxonomic lineages based on the NCBI taxonomy. Sequence functional features were assessed using several annotation tools. First, sequences were annotated with InterProscan ³¹ and InterPro³² (interproscan-5.59-91.0). Additionally, domainMapper ³³ [version 3.0.2] and the ECOD classification of known 3D structures of domains were used to complete the sequence annotation. This HMM parsing algorithm provides the detection of non-contiguous, insertional, and circularly permuted domains as well. Results from both tools were filtered with an E-value threshold of 1e-10. Transmembrane regions were annotated using deepTMHMM³⁴ [web version 1.0.20]. Finally, calcyanin sequences were detected using pCALF (see below and **Supplementary Figure 1**). Hydrophobic cluster analysis (HCA^{35,36}) was also used to assess the foldability of the analyzed sequences ^{37,38}.

Calcyanin detection and classification (Supplementary Figure 1)

Sequences of calcyanins were detected in our dataset using a dedicated in-house tool called pCALF (standing for python CALcyanin Finder). pCALF uses four hidden Markov model (HMM) profiles describing the C-terminal glycine zipper triplication specific of calcyanins, called the (GlyZip)₃ motif¹². The (GlyZip)₃ HMM profile describes the whole triplication, while Gly1, Gly2 and Gly3 HMMs describe each glycine zipper individually. These HMMs profiles were searched against amino acid sequences using pyHMMER^{39,40}. Additionally, a set of domains (7 Y-type, 1 X-type, 14 Z-type and 15 CoBaHMA), described as N-terminal domains of known calcyanins, was considered to annotate the N-terminal region of the sequences. A sequence is classified as calcyanin when it has a significative hit against the (GlyZip)3 HMM profile (sequence coverage threshold: 60%; E-value threshold: 1e-20) and significative hits against individual Gly1, Gly2 and Gly3 zippers in this specific order (sequence coverage threshold: 70%, E-value threshold: 1e-10). A sequence is also classified as calcyanin when the second glycine zipper is missing, and it contains a Y-type N-terminal domain.

Classification of the full-length sequences (Figure 2-C)

Full-length proteins were clustered using mmseqs2 41 , considering a sequence coverage threshold of 97% (coverage mode: 0, cluster mode: 0, identity threshold: 0%). This threshold was the lowest that kept the length difference between the shortest and the longest sequences in a cluster below 50 residues and ensured that the length difference between the longest and the shortest sequences of the cluster was inferior to the length of a domain. Sequences representative of each cluster, including singletons, were extracted and a self-versus-self search was performed using mmseqs2 (alignment mode: 0). Alignment results were filtered using a reciprocal coverage threshold above 70%, an E-value threshold below 1e-10 and a sequence length amplitude lower than 50 aa. The bitscore was normalized (NB) by the length of the shortest sequences comprised in the alignment 42 . Finally, alignment results were filtered with a NB threshold above 0.44 that corresponds to the median value of all NBs.

A similarity network was built with NetworkX ⁴³ using alignment results from both clustering and search with sequences as nodes and similarity as edges, weighted by NB. The network was refined iteratively in two steps: community detection and edges removal. First, the best partition was found using the Louvain Community Detection Algorithm ⁴⁴. Then, we focused on communities where the longest and shortest sequences/nodes had a length difference of more than 50 aa. We removed the node that was the furthest from the median length of the community. We repeated this process until the difference between the longest and the shortest nodes in the considered community felt under the 50 aa threshold. All the nodes removed in this way were labeled as invalid nodes. Finally, we removed all edges between invalid and valid nodes. Community detection and edges removal were repeated until no more edge could be removed.

The partitioning quality of the final network was assessed. Network layout (edge-rendering and weighted spring embedded layout, with weight interpreted as normalized values, 'strength of a disconnected spring' set to 0.01 and 'strength to apply to avoid collisions' set to 10) and rendering were produced with Cytoscape 45 .

Multiple alignment of representative sequences of the CoBaHMA family

In order to build a multiple sequence alignment (MSA) that is representative of CoBaHMA's diversity, and avoid the over-representation of a subgroup of CoBaHMA domains, the 2305 CoBaHMA sequences were clustered with mmseq2 (coverage mode: 0, cluster mode: 0, identity threshold: 60%, sequence coverage threshold: 80%)⁴¹. The representatives of the clusters and singletons were gathered, amounting to a total of 1434 CoBaHMA sequences. The sequences were aligned using mafft v7.487 ⁴⁶(maxiterate 1000; localpair), with some manual correction. The multiple sequence alignment (MSA) was viewed and analyzed with Jalview $2.11.2.6^{47}$.

WebLogo v2.8.2 ⁴⁸ was used on the Berkeley server (https://weblogo.berkeley.edu/logo.cgi) to build the logo of the MSA, restricted to the CoBaHMA β strands. In order to identify the amino acid conservation patterns, we removed the indels from alignment before building the logo.

Analysis of the 3D structure models

The AF2 3D structure models were manipulated and visualized using Chimera ⁴⁹. Based on ⁵⁰, the pLDDT values, describing the model confidence at the amino acid level, were split into 4 categories: pLDDT [?] 90 (very high confidence); pLDDT [?] [70, 90[(high confidence); pLDDT [?] [50, 70[(low confidence); and pLDDT <50 (non-interpretable). The associated predicted aligned error (PAE), extracted from AFDB ²²was also considered in order to evaluate interdomain contacts. Each position of the PAE is the uncertainty in Å on the relative positions between 2 aa, whose positions in the sequence are given by the x and y coordinates.

Structural similarities were searched using Foldseek⁵¹. Three dimensional structure models were analyzed in light of the multiple sequence alignments of each community, built using mafft v7.487 ⁴⁶ and rendered using ESPript3⁵².

RESULTS

Identification of CoBaHMA domains not contained in calcyanins

We proposed here a methodology dedicated to the specific identification of CoBaHMA domains, as members of a novel family within the large HMA domain superfamily. Starting from the sequences of CoBaHMA domains from the 15 calcyanins reported in ¹² as queries, we performed an iterative, profile-based sequence similarity search combined with the consideration of structural features provided by AlphaFold2 (AF2) models (see Material and Methods; Figure 2). Considering these structural features during the search process improved the discrimination between HMA and CoBaHMA domains. However, this specificity was achieved at the cost of a lower recovery of CoBaHMA domains, since not every sequence in the UniClust30 database had a model available in the AlphaFoldDB (AFDB) at the time of our study. We increased the size of the sequence dataset by considering transitive searches (6 iterations) and using newly detected CoBaHMA domains as probes for additional searches. A total of 38444 distinct sequences identifiers were recorded by these searches. Among these identifiers, 28918 had AF2 models. The remaining ones mostly corresponded to UniParc sequences. CoBaHMA-specific features were considered using the AF2 models, restricting the set to a total of 2358 domains corresponding to 2280 sequences, the manual inspection of which identified 68 false positive domains. Within this whole set, we thus identified a total number of 2305 (2290 domains from Uniclust30 + 15 initial probes) CoBaHMA domains within 2227 different proteins (2212 sequences from Uniclust30 + 15 initial probes).

Most of the sequences were identified during the two first iterations (see **Supplementary Figure 2** for details). Only 8 phyla are represented by more than 20 sequences. All are affiliated to Bacteria. Among them, the Proteobacteria, Cyanobacteria and Bacillota phyla are represented by 1025 (861 and 153 detected

during the two first iterations), 445 (349, 78) and 365 (274, 89) sequences, respectively. Forty-three, 87 and 40 sequences affiliated to the Actomycetota phylum were identified during the three first iterations. During the whole iterative similarity search process, only one sequence from Eukaryotes was detected during the first transitive search (Chordata, Chondrychtyes class, UniRef100_A0A401TJW7).

Figure 3 illustrates the conserved sequence patterns of the CoBaHMA domains, derived from the alignment of 1434 sequences representative of the whole domain diversity (see Material and Methods).

Hosted file

image3.emf available at https://authorea.com/users/669589/articles/669968-alphafold2-guideddescription-of-cobahma-a-novel-family-of-bacterial-domains-within-the-heavy-metalassociated-superfamily

Figure 3: Amino acid conservation in the CoBaHMA family sequences. Amino acids are colored according to their properties (black: apolar, green: polar, yellow: small, orange: aromatic, blue: basic, red: acid, cyan: histidine). The most conserved polar amino acids are displayed on the AlphaFold2 3D structure model of a CoBaHMA domain, for which per amino acid pLDDT values are mostly very high (UniProt A0A545SE61 – community C299).

The β -sheet displays several conserved features, spread over all β -strands, except strand $\beta 4$, which is not present in every CoBaHMA domain. By contrast, the two α -helices are highly variable and could not be aligned. First, strand β 1 possesses two arginines that are highly conserved: 1266, i.e. 88% of the sequences had the two arginines and 1353 (94% of the sequences) had at least one of them. These two arginines are accompanied by another basic residue (arginine or lysine) on the C-terminus of strand $\beta 1$ as well as a fairly conserved histidine (925 (65%) sequences) in strand β_0 , and together, form a basic patch at the surface of the β -sheet. The full motif HxxxRxRxR that was originally identified in calculations ¹² was present in only 27% of the CoBaHMA sequences. A continuum can thus be highlighted in the CoBaHMA family, from sequences that do not have the basic patch to sequences that have the full HXXXRxRxR motif, hinting at the possible existence of sub-families with specific features. Besides this basic patch, strand β 1 has a conserved PG motif on its N-terminus. Strand β 3 has an array of conserved small amino acids (G/A/T/S) on its N-terminus, as well as an aromatic (Y/F/H) position on its C-terminus, which is oriented toward the hydrophobic core of the CoBaHMA in the 3D structure model. Considering their nature (small or apolar), and/or their position (loops or orientations towards the hydrophobic core), all these conserved amino acids are likely of structural importance. Finally, two polar positions occupied by an asparagine (C-terminus of strand β^2) and an acidic residue (C-terminus of strand β 3) are also worth mentioning. Interestingly, strand β 2, which is the farthest from the basic patch located on strands $\beta 0$ and $\beta 1$ appears to have less sequence conservation, except a central position sometimes occupied by a basic residue. This further strengthens the hypothesis that the functional feature of the CoBaHMA is yielded by strands $\beta 0$ and $\beta 1$.

Structural and functional analysis of the full-length sequence communities

Communities of full-length sequences have been segmented based on a similarity network. The robustness of the affiliation of full-length sequences to a given community was achieved by using a refinement method described in the Material and Methods (section 2.6).

Communities can be classified into several categories, depending on whether the CoBaHMA domain alone constitutes the entire protein (what we call a "single-CoBaHMA domain protein"), or the protein containing it is longer. Moreover, in the latter case, several scenarios can be distinguished: (i) the regions of the proteins apart from the CoBaHMA domain, are already annotated by reference to profiles contained in domain databases (InterPro (IPR), Pfam); (ii) they are related to the calcyanin GlyZip motifs (pCALF tool, see Material and Methods), (iii) they cannot be annoted. An additional processing was added, by considering together communities sharing identical annotations and/or distinct annotations but corresponding to the same major functional families. We describe the communities using this analysis workflow, adding information about their 3D structures from AF2 models, as well as conserved motifs identified from multiple sequence

alignments (MSA, see **Supplementary Data 1**). We focused our analysis on the 74 large communities, containing at least 5 sequences. Altogether, they amount to a total of 1679 sequences (*i.e.* 73 % of the total number of CoBaHMA domains).

Figure 4 (see next page): Taxonomic composition of the 74 large communities of the full-length CoBaHMA domain sequences. Each bar corresponds to the taxonomic composition of a community of full-length sequences with at least 5 members. The number of sequences for a given phylum is indicated on the bar, reported to the total number of sequences per community. Communities corresponding to single CoBaHMA proteins are indicated in bold. Finally, the community C366 with one sequence belonging to Euryarcheota phylum is marked with an asterisk.

Hosted file

image4.emf available at https://authorea.com/users/669589/articles/669968-alphafold2-guideddescription-of-cobahma-a-novel-family-of-bacterial-domains-within-the-heavy-metalassociated-superfamily

Most of these 74 large communities are composed of sequences affiliated with different phyla. However, 26 of them (313 sequences) are specific to one bacterial phylum (**Figure 4**): 8, 9 and 9 communities comprised sequences affiliated to Bacillota (C675, 23 sequences; C521, 20 sequences; C4, 18 sequences; C413, 11 sequences; C603, 9 sequences; C397, 8 sequences; C369, 7 sequences; C396, 5 sequences), Proteobacteria (C201, 53 sequences; C320, 13 sequences; C562, 10 sequences; C55, 10 sequences; C584, 8 sequences; C249, 7 sequences; C344, 5 sequences; C40, 5 sequences; C7, 5 sequences), or Cyanobacteria (C283, 31 sequences; C192, 15 sequences; C20, 10 sequences; C668, 8 sequences; C706, 8 sequences; C687, 7 sequences; C628, 7 sequences; C623, 5 sequences; C654, 5 sequences) only, respectively.

While most of the 74 large communities do not have any significant functional or structural annotation (non-annotated communities), 20 InterPro (IPR) accessions and 11 ECOD accessions were detected in more than one full-length sequence within 23 large communities (**Supplementary Table 1-A**). These can be grouped into families related to P-type ATPases (10 and 7 entries from IPR and ECOD, respectively), ABC exporters (6 IPR, 2 ECOD), type 2 phosphatidylglycerol phosphatases (2 IPR, 1 ECOD) or HMA with no overlap with the CoBaHMA domains (2 IPR, 1 ECOD). A good agreement between these two sources of annotations is observed.

Moreover, it should be noted that among the 365 sparsely populated communities (composed of n<5 sequences), 154 communities have IPR or ECOD annotations shared with the larger communities (see below for details).

Overall, the 74 large communities correspond to proteins with different modular organization (Figure 5 and Supplementary Data 2 for the details in each community). These communities may be grouped into categories as follows: (i) small sequences (~100 aa length) containing one CoBaHMA domain only; larger sequences combining a CoBaHMA domain with (ii) functionally annotated regions (from IPR or other sources), (iii) non-annotated regions, or (iv) additional CoBaHMA domains (either 2 or 3). The position of the communities relative to the functional families/phyla and within the network are described inSupplementary Figure 3 and Supplementary Figure 4, respectively.

Hosted file

image5.emf available at https://authorea.com/users/669589/articles/669968-alphafold2-guideddescription-of-cobahma-a-novel-family-of-bacterial-domains-within-the-heavy-metalassociated-superfamily

Figure 5: Modular organization of the full-length CoBaHMA domain proteins. Each line illustrates the representative sequence of a large community (the vertical order follows the horizontal order in Figure 4). Line labels represent the accession of the representative sequence, the community number and the dominant phylum. Functional (IPR) and structural (ECOD) annotations are indicated along the sequence by

thin and large shaded areas, respectively. Each functional category of domain is highlighted by the following color code: CoBaHMA (green), HMA_2 (lime), HMA (yellow), P-type (cyan) and SERCA (darkblue), ABC (orange), Calcyanin Gly-Zip (red), PAP2 (magenta). Membrane regions as identified by deepTMHMM are indicated by gray areas.

Single-CoBaHMA domain proteins

Nine of the 74 large communities (C437, C201, C583, C675, C399, C397, C71, C685, C688) fall into this category, scattered over several phyla (**Figure 6**). They include a total of 422 sequences, with 90 % included in 5 communities (C437 (185 sequences), C685 (64 sequences), C201 (53 sequences), C399 (42 sequences), C583 (35 sequences)). The lengths of these sequences are around 100 aa (minimum and maximum mean lengths of 89.6 and 127.6 aa for communities C688 and C685, respectively, **Supplementary Figure 4**). Superimposition of the AF2 3D models of their representative members (in which all the core α helices and β strands are predicted with very high/high pLDDT values) indicate that these single-CoBaHMA domain proteins possess extra-regular secondary structures (mostly α helices) at the N- and/or C-terminus of the domain, which pack against the core (**Figure 6**). Some variations are observed in the β -sheet conserved motifs, with the particular case of community C397 (**Figure 6-F**), in which the conserved basic amino acids are nor present and the extra-N-terminal helix takes the place of helix α 2. Considering this topological difference, community C397 should thus be considered apart, probably not belonging to the CoBaHMA family.

Hosted file

image6.emf available at https://authorea.com/users/669589/articles/669968-alphafold2-guideddescription-of-cobahma-a-novel-family-of-bacterial-domains-within-the-heavy-metalassociated-superfamily

Figure 6: Single CoBaHMA proteins. AF2 3D structure models of the representative proteins from the nine communities of single CoBaHMA domain proteins (ribbon representations), colored according to the pLDDT values. The conserved amino acids are shown as ball-and-sticks. Core α -helices and β -strands are only labeled for the first community, with the exception of C397 (panel F), in which the extra-N-terminal helix takes the place of helix α 2. Proteins are referenced with their UniProt accession numbers: A) C437: A0A1N7C237, B) C201: A0A2D3VWG2, C) C583: A0A231P312, D) C675: W1SB97,E) C399:F5RIY2 (instead of A0A1H6HD08 (UniParc entry)),F) C397: A0A426D834, G) C71: A0A109SXH9, H)C685: A0A371IS94, I) C688: A0A0M2V0M3. The multiple sequence alignments of the communities, together with the AF2 predicted secondary structures are provided in Supplementary Data 1.

It is difficult to estimate how many sequences of single-CoBaHMA domain proteins are included in sparsely populated communities (n < 5), given that sequences of similar length (e.g. C444, mean length 123.4, **Supplementary Figure 4**) can have additional secondary structures decoupled from, instead of associated with the CoBaHMA domain, as illustrated below.

CoBaHMA domains in multidomain proteins

* IPR-annotated communities

Only a few large communities show IPR annotations. They exclusively correspond to membrane proteins. We grouped together the IPR categories relating to the same protein families (**Supplementary Table 1-A**). Worth noting, small communities (n < 5), when annotated, are mostly covered by the same IPR categories, with very few new IPR found there (**Supplementary Figure 5** and **Supplementary Table 1-B**).

1) P-type ATPases . P-type ATPases account for the majority of the total, with 345 sequences in the large communities (Supplementary Table 1-A). Below, we describe the general features of P-type ATPases and how CoBaHMA-containing P-type ATPases, forming several communities with different characteristics, are clearly distinguishable from the already well-characterized members of this superfamily.

P-type ATPases are composed of a common core of three conserved domains: (i) a discontinuous transportdomain (**T-domain**) made of six membrane-spanning helices (M1 to M6) providing the substrate translocation pathway, (ii) an ATP-binding domain (ATPBD - between M4 and M5), which includes the nucleotidebinding domain (N-domain) and the phosphorylation domain (P-domain), and (iii) an actuator domain (A-domain, between M2 and M3), which is believed to transmit changes in the ATPBD to the T-domain and to drive dephosphorvlation ¹¹. Two additional domains can complete this common core, depending on the considered P-type ATPase subset: (i) the **S-domain**, which is an auxiliary membrane unit providing support to the T domain and is located at various positions in the sequence (N- or C-terminal relative to the T domain); (ii)**regulatory (R) domains**, which are located at the N-terminus and/or C-terminus and act as intramolecular inhibitors, sensors for transported cations and/or regulators for cation affinities¹¹. Transport is accomplished via a so-called Post-Albers cycle in which phosphorylation of a conserved aspartate residue in the ATPBD causes the protein to cycle between high (E1)- and low (E2)-affinity ion-binding states. InterPro entries (IPR, Supplementary Table 1-A) are available to annotate the P-type AT-Pases over their full-length common core (IPR001757) or domains (IPR023298 : T-domain; IPR008250 : A-domain; IPR023299 : N-domain; IPR044492, IP036412, IPR023214 : HAD/HAD-like), while other IPRs provide annotations for specific P-type subsets (e.g. IPR027256 for P_{1B} -type, IPR004014 and IPR006068 for cation-transporting P-type ATPases N-terminal and C-terminal, respectively). IPR entries specific to HMA domains (IPR036163 and IPR006121) are also found, outside the limits of the CoBaHMA domains, as accompanying some P_{1B} -type ATPases

Nineteen communities with at least 5 sequences and scattered over several phyla (Figure 5) are annotated as P-type ATPases (Supplementary Table 1-A). Fifteen out of the 19 communities (C366, C429, C656, C40, C712, C525, C140, C7, C154, C692, C20, C710, C744, C550, C658) belong to the P1Btype. Indeed, the AF2 models of their representative sequences (Figure 7) include a S-domain specific to this subset, comprising two transmembrane helices, a long and curved MA helix and a kinked MB helix with an amphipathic MB' segment at the cytoplasmic membrane interface, lining the ion entry point^{53,54}. However, despite this unequivocal structural connection, only 8 out of these 15 communities (C366, C712, C658, C550, C40, C429, C525, C140) match the IPR027256 (P_{IB}-type) profile. The 15 communities differ by the architecture of their whole proteins. Five communities have a N-terminal CoBaHMA domain (illustrated in Figure 7-A with C366):C366 (114 sequences, representative member: A0A3P1Y6T0), C712 (24 sequences, representative member: A0A351Z0B2), C429 (10 sequences, representative member: A0A5C7XZV1), C40 (5 sequences, representative member: A0A7M1LI7 0) and C656 (5 sequences, representative member: A0A4P7ZPZ1). One community, C525 (61 sequences, representative member: A0A5C7KBI4), has a N terminal CoBaHMA + HMA couple (Figure 7-B). Four communities (illustrated in Figure 7-C with C140) are characterized by tandems of CoBaHMA domains forming a continuous β -sheet, with two additional β -strands in the sequence linking them: C140 (18 sequences, representative member: B8ETC0), C7 (5 sequences, representative member: A0A4Q5L4V0), C154 (5 sequences, representative member: A0A6B8M3E9) and C692 (5 sequences, representative member: A0A6B8M3E9). However, only the C-terminal CoBaHMA domain was detected by our search, while the N-terminal one lacked most of the basic residues.

Hosted file

image7.emf available at https://authorea.com/users/669589/articles/669968-alphafold2-guideddescription-of-cobahma-a-novel-family-of-bacterial-domains-within-the-heavy-metalassociated-superfamily

Figure 7: P-type ATPases. AF2 3D structure models of the representative proteins from eight P-type ATPase communities (ribbon representations), colored according to modular organization (top). These eight communities summarize the different architectures observed in large communities of P-type ATPases. Domains are designated as in¹¹. At the bottom are shown the CoBaHMA domains (left) and the M4 transmembrane helices (right), colored according to the pLDDT values and with the conserved amino acids shown as ball-and-sticks. Proteins are referenced with their UniProt accession numbers:**A)** P1B-ATPases

with a N-terminal CoBaHMA domain: C366: A0A3P1Y6T0 (for M4, the 3D structure AF2 model of D8F5K4 is also shown at right, representative of another C366 sub-community whose M4 signature sequence differs from that of subcommunity to which A0A3P1Y6T0 belongs). C366 is representative of the architecture of C429, C712, C40 and C656. **B)** P1B-ATPases with a N-terminal tandem of CoBaHMA + HMA domains: C525: A0A5C7KBI4, **C)** P1B-ATPases with a N-terminal tandem of two CoBaHMA domains: C140: A0A1Y6CVZ1. C140 is representative of the architecture of C7, C154, C692. **D)** truncated P1B-ATPases: C20: A0A3S1CNK6. C744 and C710 also belong to this group (lack of N-domain and N- and P-domains, respectively). **E-F)**P1B-ATPases with additional domains (**E:** C550: K8GFG5,**F:** C658: R5Q7W6), **G)** P1B-ATPases with unusual S-domain C413: R1CS51, **H)** Cation-transporting AT-Pases: C556: A0A7C4R520. C556 is representative of the architecture of C42 and C98. The multiple sequence alignments of the communities, together with the AF2 predicted secondary structures are provided in **Supplementary Data 1**.

Three communities lack the N and/or P domains, with the truncations matching the limits of these domains. Hence, the **C20** community (10 sequences, *representative member: A0A3S1CNK6*, illustrated in **Figure 7-D**), specific to cyanobacteria, lacks the P-domain and has degenerated consensus motifs in domains A and N, in contrast to the other communities which preserve these critical sequences (A-domain ([TS]-G-[DE]), P-domain (DKTGT) and N-domain (HP)) (**Supplementary Data 3**). **C744** (6 sequences, *representative member: A0A096BD71*, not shown) lacks the N-domain and the two C-terminal helices of the T domain, while still bearing a P-domain. **C710** (8 sequences, *representative member: A0A3D0NRW4*, not shown) lacks both the N- and P domains, as well as the two C-terminal helices of the T domain.

A remarkable point here is that of the community C369(7 sequences, representative member: A0A267MLA1), which is not detected by any P-type ATPase profile, consists of a CoBaHMa domain, the T domain and the N domain, but lacks the MA-MB-M1-M2 block, the A domain and the P domain.

Finally, two last communities correspond to P1B-ATPases including other domains in addition to an N-terminal CoBaHMA domain: C550 (5 sequences, *representative member: K8GFG5*, Figure 7-E), having a C-terminal, well predicted but yet uncharacterized domain (green) and C658 (6 sequences, *representative member: R5Q7W6*, Figure 7-F), which possesses multiple helices (however predicted with low pLDDT values) between the CoBaHMA and S-domains and after the S-domain (green).

In addition, among the 19 communities annotated as P-type ATPases, the C413 community (11 sequences, representative member: R1CS5, Figure 7-G) possesses a single N-terminal CoBaHMA domain but has an atypical MA-MB segment, which does not match the usual topology encountered in P1B-ATPases and could not be modelled accurately.

Finally, communities C556 (23 sequences, representative member: A0A7C4R520, Figure 7-H), C42 (10 sequences, representative member: A0A1M3N6C8, not shown) and C98 (13 sequences, representative member: A0A0M0SKF3, not shown) also possess a CoBaHMA domain in the N-terminus but do not belong to the P_{IB}-type. Instead, they are annotated by cation N-terminal (IPR004014) and cation C-terminal (IPR006068) profiles (Supplementary Table 1-A), which are found in several cation-transporting, P2A-ATPases (Na⁺, K⁺, Ca²⁺). Inspection of the AF2 model indicated a conserved calcium-binding site in the T-domain^{1,55}, including a conserved central glutamate. Of note is that some of the P1B-ATPases with CoBaHMA domains are detected by the cd07550 (P-type_ATPase_HM) profile from the Conserved Domain Database (CDD).

Members of the P1B-type subsets were described heretofore as specific to the translocation of heavy metal ions. They are divided into several groups based on conserved sequence motifs (in the unwound part of M4, but also in M5 and M6) and the selectivity of the transported metal ion⁵⁶. The fifteen communities of P1B-ATPases highlighted here possess conserved motifs in the unwound part of M4, which differ from one community to the other, and sometimes within large communities, in which sub-communities can be distinguished based on this feature. For instance, the biggest community **C366** can be divided into two sub-communities, according to these motifs (**Figure 7-A**): (i) Part of the sequences, such as D8F5K4, has

a conserved C-P-C motif (right in Figure 7-A), typical of heavy metal binding sites; (ii) other sequences, such as the representative sequence A0A3P1Y6T0, contain a characteristic conserved motif including an aspartate as well as a basic residue (D-[YF]-x-[TC]-x(2)-[KRH] (left in Figure 7-A). A C-P-C motif is also found conserved in community C525 , consistent with the presence of an HMA domain (in addition to the CoBaHMA domain). A large part of M4 motif within the P1B communities included a D-[YF]-x-[TC]-x(2)-[KRH] motif, with varying degrees of variability around these conserved amino acids. Only the C413 community has no strictly conserved residue in the unwound part of M4, except a central proline.

Finally, it should be noted that P-type ATPases are also abundant in sparsely populated communities (**Supplementary Table 1-B** and **Supplementary Figure 5**). This indicates a very high level of diversity that far exceeds that described based on the analysis of the more populated communities.

In conclusion, our results indicate that CoBaHMA domains form a novel family found frequently in association with P1B-ATPases, similarly to HMA domains. However, sequence signatures of heavy metal binding in M4 are not systematically found in these P1B-ATPases, giving way to other signatures, including conserved acidic and basic residues and suggesting that P1B-ATPases are not exclusively transporting heavy metals.

2) ABC exporters. CoBaHMA domains are also found in the N-terminus of type I ABC exporters. They are present in two large communities: C538 (representative member UniProt A0A2W6BRH6, 15 sequences) and C520 (representative member UniProt E3FPA8, 6 sequences).

Type I ABC exporters, formerly known as type IV exporters⁵⁷, transport a wide variety of substrates across membranes. They consist in a TMD with six transmembrane (TM) helices and a nucleotide binding domain (NBD), which form homo- or hetero-dimers, with a swapped arrangement of two TMs. Three InterPro IPR are associated with the two large communities in which type I ABC exporters are identified: IPR036640, an ABC transporter type I of the transmembrane domain superfamily: IPR027417, a P-loop containing nucleoside triphosphate hydrolase, and IPR039421, a type 1 protein exporter, encompassing both the transmembrane domain (TMD) and the nucleotide binding domain (NBD). The experimental 3D structures closest to the AF2 models of the C538 and C520 representative sequences (FoldSeek searches) are those of bacterial ABC exporters involved in the transport of various substrates, including lipid A (MsbA, pdb 7PH4), peptides (TmrAB, pdb 6RAI) or multiple drugs (Sav1866, pdb 2HYD). The CoBaHMA-containing ABC exporter sequences exhibit canonical ABC conserved motifs in the NBD (Walker A, Walker B, ABC signature), suggesting that they are active transporters (Supplementary Data 1). In contrast to other communities with CoBaHMA domains, the ABC exporter communities do not have the conserved histidine in strand $\beta 0$, while strands $\beta 0$ and $\beta 1$ include several basic amino acids (Figure 8-A). Linkers of variable length separate the N-terminal CoBaHMA domain from the TMD. These are predicted with lower pLDDT values as random coils or most often, as TMD hairpins (e.g. A0A2W6BRH6, community C538) or both (e.g. A0A6G4WXH4, community C455 or A0A7V8NHJ8, community C499), depending on the ABC exporter sequence. It is precisely the length of the linker that makes the difference between the two most populated communities, C538 and C520(Figure 8-A).

Worth noting, a group of small communities, united under the common denominator **EcsC proteins** (IPR024787: **C663** : 4 sequences, *representative member:* A0A552EVV8 ; **C684** : 1 sequence, A0A098TP70 ; **C83** : 1 sequence, A0A6H2NMM5 , **C52**: 1 sequence, Q606U9) is related to ABC transport systems. Indeed, in *Bacillus subtilis* , EcsC is found in an operon with EcsA and EcsB, which are components of an ABC transport system ⁵⁸. The AF2 model of the EcsC domain folds as an *a priori* soluble bundle of TM helices, with most of the amino acids characterized by low pLDDT values (**Figure 8-B**). A FoldSeek search did not highlight any significant similarity with known 3D structures, suggesting that this domain adopts an as yet uncharacterized fold. The EcsC CoBaHMA domains possess the characteristic His/Arg signature in the first two β strands.

In conclusion, our results indicate that CoBaHMA domains are also found in a few members of another family of membrane proteins, the ABC exporters, as well as in uncharacterized components of an ABC transport system. C) Type 2 phosphatidylglycerol phosphatases (PAP2). CoBaHMA domains are present in the Nterminus of some type 2 phosphatidylglycerol phosphatases (PAP2) in Cyanobacteria. Communities C419 (representative member UniProt A0A1U7ILD7) and C643 (representative member UniProt A0A856MGD8) contain 24 and 6 sequences, respectively. PAP2 sequences are described by two InterPro entries: IPR036938 and IPR000326, both entries being defined as phosphatidic acid phosphatase type 2/haloperoxidase. Integral membrane proteins from the PAP2 family dephosphorylate a variety of compounds, including lipids and carbohydrates ⁵⁹. They consist in a core TMD, with six tightly packed TM helices connected by extramembrane loops, two of which interacting together to form the catalytic site (Figure 8-C). The sequences of the CoBaHMA-containing PAP2 belong to the lipid phosphatase/phosphotransferase (LPT) family, as they all contain the conserved tripartite active site motif $(KX_6RP-PSGH-SRX_5HX_3D)^{60}$ (Supplementary Data 3). Members of this family modify several types of lipids in Gram-negative bacteria, e.q. phosphatidylglycero-phosphate (PGP) for PgpB⁶¹, or lipid A for LpxE⁶². The AF2 model of the representative sequence resembles the 3D structure of B. subtilis bsPgpB (pdb 6FMX(A); Prob 1.00, 26.9 % identity), which contains eight TM α -helices, six of them (α 1, α 4- α 8) being tightly packed, while the α^2 helix is amphiphilic, lying at the surface of the lipid bilayer on the active site side (Figure 8-C). The linker separating the CoBaHMA domain, located in the intracellular milieu, from the terminal TM α 1-helix is variable. It is predicted as a random coil or even as integrating additional TM α -helices (e.g. A0A433NH73, community C542), always with very low pLDDT values. Again, the PAP2 CoBaHMA domains possess the characteristic His/Arg signature in the first two β strands.

Hosted file

image8.emf available at https://authorea.com/users/669589/articles/669968-alphafold2-guideddescription-of-cobahma-a-novel-family-of-bacterial-domains-within-the-heavy-metalassociated-superfamily

Figure 8: ABC exporters - EcsC proteins - PAP2 - DAGK . AF2 3D structure models of the representative proteins from communities including ABC exporters, EcsC proteins, type 2 phosphatidyl-glycerol phosphatases (PAP2) and diacylglycerolkinases (DAGK), colored according to the pLDDT values. CoBaHMA domains are shown at the bottom, with conserved amino acids shown as ball-and-sticks. Proteins are referenced with their UniProt accession numbers: A) C538: A0A2W6BRH6, B) C663: A0A552EVV8, C) C419: A0A1U7ILD7, D) C515: A0A0S3UCN3,E) C553: A0A841V906. The 3D structures of the ABC exporter dimer model (A) was built after superimposition of the AF2 model single chain on the experimental 3D structure of TM287/TM288 (best hit in a HH-PRED search, respecting the distance between the two swapped TMs and the TMD core). The multiple sequence alignments of the large communities, together with the AF2 predicted secondary structures are provided in **Supplementary Data 1**.

D) Diacylglycerol kinases

Although we mostly focused on large communities, two small communities of cyanobacterial proteins, C515 (1 sequence, A0A0S3UCN3) andC553 (2 sequences, representative member: A0A841V906), caught our attention. They match several IPR profiles (IPR045540; IPR017438; IPR016064; IPR005218; IPR001206; IPR004363) related with bacterial diacylglycerol (DAG) kinases. These enzymes convert DAG, formed by the turnover of membrane phospholipids, to phosphatidic acid ⁶³. In these proteins, the CoBaHMA domain is located at an unusual C-terminal position (Figure 8-D), while the N-terminal domain corresponds to the catalytic DAK kinase (FoldSeek matches with the putative DAG kinase from *Bacillus anthracis* (pdb 3T5P(B), Prob. 1.00, 25.2 % identity, Hou et al. unpublished) and the DAG kinase DgkB from Staphylococcus aureus (pdb 2QVL(A), Prob. 1.00, 22 % identity ⁶³). The 3D structure of DgkB has a two-domain architecture, similar to that found in *E. coli* YegS, which phosphorylates *in vitro* phosphatidyl glycerol⁶⁴. Members of the C515 and C553families share the conserved P-loop (φ -x-x-G-G-D-G-T- φ , where φ represents a hydrophobic amino acid), but exhibit slight differences in the two other conserved motifs, as described in ⁶³((φ -ph-x-N-P-x-S/A-G instead of φ -ph-x-N-P-x-G-T-x-N-A- φ -x-N instead of φ -ph-P-x-G-T-x-N-D- φ -x-R; side and top of the nucleotide-binding site, respectively). All three sequences share a common domain in between the DAG kinase and the CoBaHMA domains, modeled with very low

pLDDT values as a long helix (A0A841V906) or a two-helix hairpin (A0A0S3UNC3). The C515 single member community (A0A0S3UCN3) differs from the two other sequences of the C553 community by an additional C-terminal domain, which is related to the GlyZip motif described below. The DAGK CoBaHMA domains possess the characteristic His/Arg signature in the first two β -strands.

* Other annotated, non-IPR, communities

A) HMA_2

Nine hundred and sixty-three CoBaHMA domains overlap the Pfam profile HMA_2 (PF19991). Some communities are almost entirely covered by the profile, while others are only partially. This profile is described as distantly related to HMA domains in its N-terminal part, containing in particular the conserved histidine we also highlighted in this study. Not all the CoBaHMA domains match the N-terminal part of the HMA₂ profile, indicating that proteins matching this profile constitute a subset of the CoBaHMA family. However, the HMA_2 profile is larger than CoBaHMA domains, with a total length of 180-190 aa, and includes at its C-terminal part a conserved region generally predicted by AlphaFold2 as two contiguous helices. Matches to this C-terminal region are observed for sequences in 13 large communities (C75, 80 sequences, representative member: A0A564ZMZ6; C393, 62 sequences, representative member: A0A1Z4FYZ0 ; C141, 46 sequences, representative member: A0A1J1CTN7; C586, 24 sequences, representative member: C9KN15; C329, 23 sequences, representative member: A0A7Y6UHT; C521, 20 sequences, representative member: A0A366XQZ5; C30, 17 sequences, representative member: A0A2U3KUJ9; C562, 10 sequences, representative member: A0A4R3M3I3; C603, 9 sequences, representative member: A0A1G9U165; C702, 8 sequences, representative member: F5RIY1; C294, 7 sequences, representative member: A0A2V7B4L6, C344, 5 sequences, representative member: A0A4P2PVK5; C5, 5 sequences, representative member: A0A662ZLY5). For communities having a C-terminal region associated with high AF2 pLDDT values ([?]70), such as C393, the two helices, often predicted as transmembrane segments, pack together to form a hairpin (Figure 9-A). However, no obvious similarity with any known 3D structure could be detected by FoldSeek for this case, outside helix hairpins belonging to larger assemblies.

Hosted file

image9.emf available at https://authorea.com/users/669589/articles/669968-alphafold2-guideddescription-of-cobahma-a-novel-family-of-bacterial-domains-within-the-heavy-metalassociated-superfamily

Figure 9: CoBaHMA domains associated with helical segments, including helical hairpins. AF2 3D structure models of the representative proteins from communities of CoBaHMA domains associated with helical segments, including hairpins (ribbon representations), colored according to the pLDDT values. CoBaHMA domains are shown in the same orientation. The communities are grouped according to their matching with the HMA_2 N-ter/C-ter profile (Pfam 19991) and the specificities of the segment accompanying the CoBaHMA domain. On the right, the AF2 models of the representative members of communities indicated in bold are shown. Proteins are referenced with their UniProt accession numbers: A) C393: A0A1Z4FYZ0, B) C202: A0A5C7T941, C) C331: Q8YVH2 (C628* lacks the central helical segment present in the three other communities), D) C341: A0A6M0J010, E) C352: A0A7Y4FMF1, F) C636: A0A1M4UKF2,G) C487: H8GQW7G, H) C740: A0A4R3PQS8, I)C192: Q8DMV2. The multiple sequence alignments of the communities, together with the AF2 predicted secondary structures is provided in

Supplementary Data 1.

Matches with the HMA_2 profile were also detected in other communities, however limited to the N-terminal CoBaHMA domain. A few of them correspond to communities including single CoBaHMA domain proteins (C339, C583, C685, C675) and P-type ATPases (C140, C369, C413) (see before). We analyzed the regions C-terminal to the HMA_2 N-ter/CoBaHMA domain of the remaining large communities (C4, C55, C202, C283, C320 C331, C341, C444, C352, C583, C584, C623, C628, C636, C654, C668, C687, C706), especially

for detecting possible distant relationships to the HMA_2-Cter profile. Three communities **C202** (57 sequences, representative member: A0A5C7T941), **C320** (13 sequences, representative member: A0A0X8JQ54) and **C584** (8 sequences, representative member: A0A3B9QA35) also possess a hairpin of two helices, often predicted as transmembrane segments, with conserved amino acids. However, they differ from those defining the HMA_2-Cter profile (illustrated with the representative sequence of the **C202** community on **Figure 9-B**). In particular, the C202, C320 and C584 C-terminal regions are characterized by highly conserved histidine residues, together with basic (arginine, lysine) residues. Moreover, a FoldSeek search against PDB detected significant similarities (Prob 0.97, 19 % identity) between the C202 representative sequence and the long alpha-hairpin domain of a manganese/iron superoxide dismutase (pdb 4BR6), which forms a four-helix bundle through protein dimerization and provides histidine residues to the ion-binding site at the interface with the preceding domain ⁶⁵.

Communities C331 (17 sequences, representative member: Q8YVH2), C668 (8 sequences, representative member: A0A0C2QKB5) and C687 (7 sequences, representative member: A0ZAK5) possess two HMA_2 N-ter/CoBaHMA domains. The first of these domains is also followed by a helical hairpin, predicted with low pLDDT values. This helical hairpin contains conserved charged and aromatic amino acids. The second CoBaHMA domain is also followed by a helical region, although less defined (Figure 9-C). In community C628 (7 sequences, representative member: A0A1Z4S7A2), two HMA_2 N-ter/CoBaHMA domains are also present, however apparently without the central helical hairpin.

A few communities are characterized by a HMA_2 N-ter/CoBaHMA domain, followed by one or several non-packed α -helices predicted with low pLDDT values : C341 (23 sequences, representative member: A0A6M0J010 - three non-packed α -helices, containing conserved basic and acidic residues (Figure 9-D), C444 (101 sequences, representative member: A0A522V264 - two non-assembled a-helices, with less amino acid conservation), C55 (10 sequences, representative member: I3YG68 – one helical segment with conserved basic and aromatic amino acids).

Representative sequences from communities C352 (31 sequences, representative member: A0A7Y4FMF1), C725 (12 sequences, representative member: A0A6F9WUC5) and C706 (8 sequences, representative member: A0A1Q4RV65) are apparently disordered (Figure 9-E). However, the two first ones contain conserved charged residues accompanying hydrophobic clusters, suggesting a hidden fold (Supplementary Figure 6).

Interestingly, two communities: C283 (31 sequences, *representative member:* B1WT30) and C636 (24 sequences, *representative member:* A0A1M4UKF2) have a helical hairpin, with conserved charged/polar residues at their N-terminus, upstream the HMA_2-Nter/CoBaHMA domain (Figure 9-F).

The remaining communities possessing HMA_2-Nter/CoBaHMA domains followed by C-terminal more complex architectures (C4, C93, C192, C623) are described in the two next chapters. Finally, we extended our analysis (in search of hairpin-like motifs) within CoBaHMA -containing sequences that do not match the HMA_2 profile. Members of the community C678 (13 sequences, representative member: A0A0M0SKF3) also contain a helical hairpin with two conserved histidine residues in the second helix (Figure 9-G). These features are similar to the ones observed in the C320 (13 sequences, representative member: A0A0X8JQ54) and C584 (8 sequences, representative member: A0A3B9QA35) communities. A helical hairpin, with highly conserved charged amino acids, is also present in three other non-HMA_2 communities: C106 (9 sequences, representative member: A0A7C2VAZ9), C578 (7 sequences, representative member: A0A1H4BM11) and C487 (47 sequences, representative member: H8GQW7). Finally, non HMA_2 CoBaHMA domains are also found associated in other helical segments -Figure 9-H):C396 with one helix (5 sequences, representative member: A0A3D2YH28), and C740 with four helices, including a hairpin (10 sequences, representative member: A0A4R3PQS8).

In conclusion, our results indicate the presence of a large number of sequences with a region downstream of the CoBaHMA domain predicted to possess two helices, which have a strong propensity to form a hairpin and some of which possess charged residues.

B) Calcyanins: a (GlyZip)3 motif detected by a dedicated tool (pCALF)

Community C192 (15 sequences, *representative member: Q8DMV2*) corresponds to calcyanins, in which the CoBaHMA domain was first identified. Some of the CoBaHMA domains of this community match the HMA_2 N-ter profile described above. Calcyanins contain a C-terminal domain, consisting in a three-fold repeat of a large glycine-zipper (GlyZip) motif. This large GlyZip motif itself corresponds to a duplicated smaller glycine zipper motif, interrupted in its middle part by a conserved Gly-Pro dipeptide. AF2 modeled this GlyZip motif as a hairpin of tightly packed helices, consistent with the presence of conserved glycine residues repeated every four residues (**Figure 9-I**)^{14,66}. The low/very low pLDDT values may be due to the very large sequence distance between this GlyZip motif and known hairpins of this type, present in different architectures (as explored with FoldSeek). However, AF2 fails to assemble these GlyZip motifs in a consistent way.

* Other (HMA₂ and non-HMA2) N-ter communities

Last, very few large communities other than those described above have been detected. Protein segments associated with CoBaHMA domains correspond to disordered or ordered regions. The order (foldability) in these regions was assessed by examining AF2 predictions, in particular the segments associated with high/very high pLDDT values ([?]0.7), as done for the HMA_2 N-ter containing proteins. By this way, we also retrieved helical hairpins in communities other than those containing HMA_2 CoBaHMA (**Figure 9-G**). In contrast, low values of pLDDT are sometimes indicative of disorder, although in some cases, low pLDDT values are associated with genuinely well folded regions (predicted as folded or in random coil), but for which the prediction cannot be supported. For example, this is the case of new folds and sequences lacking homologs ^{37,38}. One way to evidence these" hidden" folded domains is to assess foldability using Hydrophobic Cluster Analysis (HCA) ^{37,38}. Therefore, we investigated unannotated sequences by combining AF2 structure predictions with HCA analyses of the protein sequences, when needed.

The most populated (non-HMA_2 N-ter) community (C233 ; 39 sequences; representative sequence: U2SKL7) is predicted by AlphaFold2 as a CoBaHMA domain with a C-terminal extension comprising three strands completing the core β -sheet (Figure 10-A). A four-helix bundle predicted with high pLDDT values is inserted in between the CoBaHMA domain and this C-terminal extension. Interestingly, a FoldSeek search of this four-helix bundle, which contains strictly conserved histidine, arginine and aromatic residues, indicated a possible structural relationship with the MA-MB-M1-M2 block of P1B-ATPases (pdb 4UMW, Prob 0.30, 11.2 % identity, Figure 10-A). The MB kink (resulting in MB') is not present in the C233 four-helix bundle, whereas the basic residues (histidine and arginine) are located at the entry site of the P1B-ATPase.

Hosted file

image10.emf available at https://authorea.com/users/669589/articles/669968-alphafold2guided-description-of-cobahma-a-novel-family-of-bacterial-domains-within-the-heavymetal-associated-superfamily

Figure 10: CoBaHMA domains with unknown regions. AF2 3D structure models of the representative proteins from communities of CoBaHMA domains associated with unknown regions (ribbon representations), colored according to the pLDDT values. Proteins are referenced with their UniProt accession numbers: A) C233: U2SKL7, the CoBaHMA domain is followed by an helical bundle, which superimposed with the MA/MB/M1/M2 bundle of P1B-ATPases (left), B) C4: A0A1H0C0W7, the CoBaHMA domain is followed by an helical bundle, with no striking similarities with any known 3D structures, C) C93: A0A552LCK1, the CoBaHMA domain is preceded by an all-alpha domain, with no striking similarities with any known 3D structures, D) C654: A0A0M0SSI8, three CoBaHMA domains are separated from each other by disordered linkers. The two first domains (italics) lack the basic signature in strands $\beta 0$ and $\beta 1$. E) C623: A0A1Z4TPU0, the CoBaHMA domain is followed by a domain belonging to the BPI family (TULIP superfamily). The AF2 3D structure model is compared after superimposition to the experimental 3D structure of human BPI with two bound phosphatidylcholine (pdb 1BP1). The multiple sequence alignments of the communities, together with the AF2 predicted secondary structures is provided in **Supplementary Data**

1.

The C4 HMA_2 N-ter community (18 sequences, representative member: A0A1H0C0W7) is also predicted by AlphaFold2 as composed of a N-terminal CoBaHMA domain and a four-helix bundle (Figure 10-B), also with conserved amino acids (especially arginine and histidine residues) located at the tip of the bundle. However, in this case, no significant structural similarity was found between the representative member of the community and any experimental 3D structures using FoldSeek.

A third community (C93, 6 sequences, *representative member:* A0A552LCK1) possesses a C-terminal, HMA_2 N-ter CoBaHMA domain, preceded by a globular all- α domain, with conserved residues, but also which does not share obvious similarity with any available experimental 3D structure (FoldSeek search) (Figure 10-C).

Repeated CoBaHMA domains within a single protein sequence seems to be a relatively frequent feature of the family, since we also observed in a fourth community, C654 (5 sequences, *representative sequence:* A0A0M0SSI8), a C-terminal HMA_2 N-ter/CoBaHMA domain (Figure 10-D), preceded by a tandem of CoBaHMA-like domains (devoid of the conserved basic signature). However, it did not contain any other folded domain based on AlphaFold2 modelling and HCA analysis.

Finally, a fifth, community (C623, 5 sequences, representative member: A0A1Z4TPU0) with a HMA_2 N-ter, and present in cyanobacteria, possesses a C-terminal well-folded domain, which shares a clear structural relationship with the BPI (Bacterial Permeability-increasing)-like family, as indicated by FoldSeek (e.g. A0A1Z4TPU0 - pdb 1BP1 (BPI, 67), Prob. 0.98, 7.5 % identity, Figure 10-E , B7KGC3 – pdb 2OBD (CETP, Qiu et al. 2017), Prob. 1, 10.2 % identity). Members of the BPI-like family share a common fold, consisting in a long α -helix wrapped in a highly curved sheet, made of long, antiparallel β -strands and display a tubular cavity, in which lipids bind 68,69 .

Discussion

The large-scale predictions of 3D structures now enabled by AI-based approaches allow to functionally annotate large sets of proteomes at the amino acid level, and identify new folds (e.g. ⁷⁰). The 3D models provided by AlphaFold2 (AF2) via a dedicated database²², connected to UniProt²⁹, offer an unprecedented tool for extending the exploration of the universe of protein domains and studying their evolutionary trajectory. Here, we have taken advantage of this large-scale structural information, combined with sequence similarity search and clustering tools, to identify the members of a new family of domains called CoBaHMA. This domain family shares a common evolutionary origin with HMA domains, as evidenced by the signature of a common hydrophobic core. However, CoBaHMA can be discriminated from HMA based on an additional, external β strand (called β 0) and, in many cases, a specific sequence signature, conferring a positive electrostatic charge on one side of the beta sheet, which is likely to be associated with a specific function. The AF2 models provided a structural constraint throughout the workflow we built, enabling a fine discrimination between CoBaHMA domains and the rest of the very large HMA superfamily. It should be noted, however, that the structural features we have automatically defined lightly suffer from the definition of cutoffs. As a result, it is possible that highly divergent members of the CoBaHMA with large loops between the two first β strands may be overlooked. The methodology developed is also dependent on the availability and accuracy of AF2 models. Finally, the proteins identified are extracted from a non-redundant set of sequences (UniRef30). limited to sequences with no more than 30% identity. Therefore, the CoBaHMA family described in this work constitutes a minimum set, limited to these non-redundant sequences and excluding members for which no AF2 models are available and which are highly divergent.

Besides providing information about a novel family of domains, our study illustrates how evolution may operate within a superfold to provide broad functional diversity. Nevertheless, the function of the CoBaHMA family of domains has yet to be defined. Further prediction of this CoBaHMA-specific function, or at least the biological environment in which it is performed, can be aided by analyzing the domain architecture of the proteins containing the CoBaHMA domain. Here, deciphering this architecture was again aided by the AF2 predictions, combined with domain database (InterPro) annotations. A large part of the CoBaHMA com-

without permission. — https://doi.org/10.22541/au.169596218.88011536/v1 — This a preprint and has not been peer reviewed. Data may be prelin

munities corresponds to single CoBaHMA domain proteins. This is reminiscent of the single-HMA domain proteins, behaving as chaperones⁷¹. CoBaHMA domains are also found in association with a limited number of protein families, which mostly correspond to membrane proteins, at least for those annotated through InterPro profiles and predictors of transmembrane segments. These families were especially analyzed in the large communities we have described here, and they also represent the bulk of the sparsely populated communities (after the analysis of the InterPro annotations, **Supplementary Table 1-B** and **Supplementary Figure 5** for details). Only a few other IPRs in addition to the previously described EcsC and DAGK are found in these last communities, generally limited to singletons. Of note is that CoBaHMA domains are principally located at the N-terminus in multidomain proteins. CoBaHMA domains appear to be specific to bacteria, in contrast to HMA domains which are found in bacteria, archaea and eukaryotes. Overall, these observations suggest a membrane-related function of the CoBaHMA domains, specific to bacteria. Moreover, CoBaHMA domains show lineage-specific expansions, as communities are restricted to certain species or phyla. This suggests that they can accommodate unique function(s), linked to specific environments.

Like HMA domains, CoBaHMA domains are particularly abundant in P1B-ATPases. The latter are commonly defined as integral membrane proteins that couple ATP hydrolysis to the transport of metal cations, such as copper, zinc and cobalt ⁵⁶. The specificity of P1B-ATPases towards heavy metals is linked to conserved motifs in the middle of the fourth transmembrane helix (TM4). These motifs directly coordinate the ion through cysteine/histidine side chains. Besides, soluble N- and C-terminal metal-binding extensions (known under the generic term Heavy Metal Binding Domains (HMBDs)), also rich in cysteine/histidine and including HMA domains, seemingly play a regulatory role ^{53,54,72}. In particular, HMBDs interact with the amphipathic helix MB', lying at the membrane-cytosol interface at the end of a P_{IB}-specific MA and MB membrane hairpin ⁷³. This amphipathic MB' helix is connected to the high affinity ion-binding site through a conserved electronegative funnel. HMBDs also interact with the cytosolic domains (A and P domains), thus playing a potential regulatory role by interfering with conformational changes coupling ATP hydrolysis with ion transport across the membrane ⁷³. Here, we show that P1B-ATPases are not restricted to heavy metals. Indeed, the proteins identified here share a P1B-specific MA and MB membrane hairpin but contain CoBaHMA domains instead of the usually encountered HMA domains and have conserved motifs in TM4 different from those coordinating heavy metals. These motifs vary depending on the considered community, often including charged (acidic and/or basic) or polar (serine/threonine) residues. Interestingly, a few communities contain the TM4 CPC motif, typical of heavy metal transport, together with a CoBaHMA domain only, and no HMA domain at its N-terminus. This suggests that the coupling of HMA with TM4 Cysrich motif is not necessary as initially thought. Moreover, a large sequence/structure diversity is observed at the level of the MA/MB hairpin, probably linked with the substrate diversity of CoBaHMA-containing P1B-ATPases, as also evidenced by the diversity of the M4 conserved motifs. Finally, it is worth noting that CoBaHMA domains are also found associated with members of the P2A-ATPase family (SERCA, Ca²⁺ ATPases), and are thus not restricted to the P1B subgroup.

The question remains as to what are the substrate specificities of these P-type proteins associated with CoBaHMA domains, and what is the role of the CoBaHMA domains in this specific modular organization. A regulatory role similar to that played by HMA could be expected, supported by the fact that some AF2 models displayed significant interfaces with A domains, involving amino acids outside the conserved basic patch (**Supplementary Figure 7**). The identity of the ligands for these charged amino acids on the CoBaHMA domain surface of P-type proteins remains to be explored.

Accessory domains provide an additional, often regulatory effect on the functions of ATP-binding cassette (ABC) exporter cores, which are formed by transmembrane domains (TMDs) and nucleotide-binding domains (NBDs)⁷⁴. For instance, the cytosolic Cystathionin Beta Synthase (CBS) domains, at the C-terminus of the osmoregulatory OpuA, inhibits the transporter activity by binding to cyclic-di-AMP⁷⁵. This protein is gated by ionic strength, which modulates the interaction of positively charged amino acids in the NBDs with negatively charged lipids ⁷⁵. Although, similar to what we described for the P-type ATPases, the specific function of CoBaHMA domains in this ABC context remains to be discovered, a number of points can be considered. First, the ABC transporters with the highest sequence identities (~30 %) include lipid

transporters such as MsbA, suggesting that (i) the CoBaHMA-domain-containing ABC transporters might be involved in lipid transport, and (ii) the CoBaHMA domain might be directly involved in the uptake of lipids. Second, the specific position of the domain, N-terminal to the NBD, places it at the right location to interact with the polar heads of lipids, as observed for instance with the lasso domain found in some ABCC transporters such as the Cystic Fibrosis Transmembrane conductance Regulator (CFTR) protein (ABCC7) (⁷⁶ for a review). This suggests a specific role in contacting the membrane via the conserved basic patch at the surface of the domain. The predicted presence of additional TM helices between the CoBaHMA domain and the TMD in most of the ABC communities, like in the ABCC transporters ABCC1 (MRP1) and ABCC8 (SUR1) ⁵⁷, might serve an additional regulatory purpose. This "lipid hypothesis" is also interesting to consider with regard to the function of the CoBaHMA domain in the context of P1B-ATPases, particularly as specific transport of lipids is carried out by another class of P-type ATPases (P₄-type, ^{77,78}).

An interesting domain architecture is observed in Community C623, which includes a domain belonging to the BPI (bactericidal/permeability-increasing protein)-like family. This family, comprising lipopolysaccharidebinding protein (LBP) and the lipid transfer proteins CETP and PLTP, includes one or two tandem copies of a fold providing a tubular cavity for the binding of lipids⁶⁸. It was extended to more divergent groups of proteins, including the SMP (synaptotagmin-like, mitochondrial and lipid-binding proteins) domains, which are associated with eukaryotic membrane processes ⁶⁹. All the so-described families were grouped into a single superfamily called TULIP, after TUbular LIpid binding Protein) domain ^{69,79,80}. Cyanobacterial proteins of community C623 bear only one copy of a BPI-like domain, with the CoBaHMA domain likely positioned nearby a potential lipid-binding site.

The hypothesis of the CoBaHMA domain serving as a binding module for positively charged lipids (phospholipids) in bacteria is appealing considering the wide knowledge about phospholipids-binding domains in eukaryotes. Indeed, the tray of basic amino acids offered by CoBaHMA domains resembles that observed in C2 domains for instance, which interact with membranes in a Ca^{2+} -dependent manner through a polybasic cluster, with specificity to phosphatidylinositol-4,5-bisphosphate ^{81,82}.

The phosphatidylinositol phosphates that are largely recognized in eukaryotic membranes are also the targets of some bacterial proteins acting as effectors or toxins (e.g. the lipid raft targeting domain of the Bordetella pathogens ⁸³, also see⁸⁴ for a review). However, bacterial membranes are distinct in lipid composition from eukaryotic membranes, and their lipid-binding modules are far less well known. Phosphatidylglycerol (PG) might be a target for the CoBaHMA domain. This phospholipid is present in both Gram - and Gram + bacteria, and plays a central role in the synthesis of cardiolipin (CL, diphosphatidylglycerol), lysophosphatidylglycerol (LPG) and oligosaccharides⁸⁵. Phosphatidic acid (PA) is another potential candidate. It is linked to the activity of diacylglycerol (DAG) kinase, and serves as a precursor for glycerolipids ⁸⁶. An appealing hypothesis is that, in addition to being specifically recognized by the CoBaHMA domains, these phospholipids could be transported by membrane systems in which CoBaHMA is included (e.g. ABC exporter), a mechanism which could contribute to the general lipid homeostasis. However, such hypotheses remain highly speculative and need to be extensively tested.

Phosphoglycerolipids are far less abundant in cyanobacterial membranes, with PG being the only phospholipid present ⁸⁷. It is present in the thylakoid membranes in low proportion (10 %) relative to the more abundant glycerolipids monogalactosyldiacylglycerol (MGDG), digalactosyldiacylglycerol (DGDG) and sulfoquinovosyldiacylglycerol (SQDG) ⁸⁷. PG proportion is regulated in response to phosphate availability ⁸⁸, but is essential not only for photosynthesis ⁸⁹ but also cell division and metabolism ⁹⁰. We note that CoBaHMA domains are found in cyanobacteria-restricted communities of phosphatidic acid phosphatases (PAP2) and diacylglycerol (DAG) kinases. Both types of enzymes are involved in the biosynthesis of lipids starting from phosphatidic acid (PA).

One of the outstanding features of our domain grammar analysis was the co-occurrence of the CoBaHMA domain not only with already annotated, membrane-specific domains, but also with hairpins of two consecutive helices. These helical hairpins show both a conserved structural motif as revealed by AF2 models, and a wide variety of sequences: some include a lot of strong hydrophobic amino acids, as in the HMA_2 C-ter profile, and are predicted as forming transmembrane segments; others include small but also globally apolar residues. However, most of them include conserved, charged residues, in particular histidine and basic residues. This suggests that these hairpins may serve as platforms on which amino acids can be grafted to interact with specific ions or ligands. From an evolutionary point of view, it is tempting to speculate that these hairpins can be used as basic units for integrating more complex architectures, such as ABC exporter TMDs or those present in calcyanins (three repeats of a glycine-zipper helical hairpin). These calcyanin glycine-zippers are structures characterized by very compact assemblies due to the presence of glycine every 4 residues, but it remains yet to be specified whether they are soluble or membrane-bound. One open question is to what extent the MA-MB hairpin specific to P1B-ATPases may have evolved from this basic module. Finally, from a methodological point of view, it would be interesting to consider this CoBaHMA-specific grammar to improve sequence similarity searches, as done for instance by Terrapon et al. ⁹¹ and Faure and Callebaut¹⁸ or, more recently, by Buchan and Jones⁹² using natural language word embedding techniques.

References

1. Orengo CA, Jones DT, Thornton JM. Protein superfamilies and domain superfolds. *Nature*. 1994;372(6507):631-634.

2. Chitturi B, Shi S, Kinch LN, Grishin NV. Compact Structure Patterns in Proteins. J Mol Biol. 2016;428(21):4392-4412.

3. Kolodny R. Searching protein space for ancient sub-domain segments. *Curr Opin Struct Biol.* 2021;68:105-112.

4. Chandonia JM, Guan L, Lin S, Yu C, Fox NK, Brenner SE. SCOPe: improvements to the structural classification of proteins - extended database to facilitate variant interpretation and machine learning. *Nucleic Acids Res.* 2022;50(D1):D553-d559.

5. Caetano-Anollés G, Caetano-Anollés D. An evolutionarily structured universe of protein architecture. Genome Res.2003;13(7):1563-1571.

6. Thornton JM, Orengo CA, Todd AE, Pearl FM. Protein folds, functions and evolution. J Mol Biol. 1999;293(2):333-342.

7. Grishin NV. Fold change in evolution of protein structures. J Struct Biol. 2001;134(2-3):167-185.

8. Jung J, Lee B. Circularly permuted proteins in the protein structure database. *Protein Sci.* 2001;10(9):1881-1886.

9. Arnesano F, Banci L, Bertini I, et al. Metallochaperones and metal-transporting ATPases: a comparative analysis of sequences and structures. *Genome Res.* 2002;12(2):255-271.

10. Bull PC, Cox DW. Wilson disease and Menkes disease: new handles on heavy-metal transport. *Trends Genet.* 1994;10(7):246-252.

11. Palmgren MG, Nissen P. P-type ATPases. Annu Rev Biophys.2011;40:243-266.

12. Benzerara K, Duprat E, Bitard-Feildel T, et al. A New Gene Family Diagnostic for Intracellular Biomineralization of Amorphous Ca Carbonates by Cyanobacteria. *Genome Biol Evol.* 2022;14(3).

13. Kim S, Chamberlain AK, Bowie JU. Membrane channel structure of Helicobacter pylori vacuolating toxin: role of multiple GXXXG motifs in cylindrical channels. *Proc Natl Acad Sci U S A*.2004;101(16):5988-5991.

14. Kim S, Jeon TJ, Oberai A, Yang D, Schmidt JJ, Bowie JU. Transmembrane glycine zippers: physiological and pathological roles in membrane proteins. *Proc Natl Acad Sci U S A*.2005;102(40):14278-14283.

15. De la Concepcion JC, Franceschetti M, Maqbool A, et al. Polymorphic residues in rice NLRs expand binding and response to effectors of the blast pathogen. *Nat Plants.* 2018;4(8):576-585.

16. Jiang D, Zhao Y, Fan J, et al. Atomic resolution structure of the E. coli YajR transporter YAM domain. *Biochem Biophys Res Commun.*2014;450(2):929-935.

17. Jiang D, Zhao Y, Wang X, et al. Structure of the YajR transporter suggests a transport mechanism based on the conserved motif A. Proc Natl Acad Sci U S A. 2013;110(36):14664-14669.

18. Faure G, Callebaut I. Identification of hidden relationships from the coupling of hydrophobic cluster analysis and domain architecture information. *Bioinformatics*. 2013;29(14):1726-1733.

19. Jin J, Xie X, Chen C, et al. Eukaryotic protein domains as functional units of cellular evolution. *Sci* Signal.2009;2(98):ra76.

20. Vogel C, Bashton M, Kerrison ND, Chothia C, Teichmann SA. Structure, function and evolution of multidomain proteins. *Curr Opin Struct Biol.* 2004;14(2):208-216.

21. Jumper J, Evans R, Pritzel A, et al. Highly accurate protein structure prediction with AlphaFold. *Nature*.2021;596(7873):583-589.

22. Varadi M, Anyango S, Deshpande M, et al. AlphaFold Protein Structure Database: massively expanding the structural coverage of protein-sequence space with high-accuracy models. *Nucleic Acids Res.* 2022;50(D1):D439-d444.

23. Mirdita M, von den Driesch L, Galiez C, Martin MJ, Söding J, Steinegger M. Uniclust databases of clustered and deeply annotated protein sequences and alignments. *Nucleic Acids Res*.2017;45(D1):D170-d176.

24. Steinegger M, Meier M, Mirdita M, Vöhringer H, Haunsberger SJ, Söding J. HH-suite3 for fast remote homology detection and deep protein annotation. *BMC Bioinformatics*. 2019;20(1):473.

25. Kabsch W, Sander C. Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers*. 1983;22(12):2577-2637.

26. Touw WG, Baakman C, Black J, et al. A series of PDB-related databanks for everyday needs. *Nucleic Acids Res.* 2015;43(Database issue):D364-368.

27. Rost B, Eyrich VA. EVA: large-scale analysis of secondary structure prediction. *Proteins.* 2001;Suppl 5:192-199.

28. Ettema TJ, Huynen MA, de Vos WM, van der Oost J. TRASH: a novel metal-binding domain predicted to be involved in heavy-metal sensing, trafficking and resistance. *Trends Biochem Sci.*2003;28(4):170-173.

29. UniProt: the Universal Protein Knowledgebase in 2023. Nucleic Acids Res. 2023;51(D1):D523-d531.

30. Shen W, Ren H. TaxonKit: A practical and efficient NCBI taxonomy toolkit. J Genet Genomics. 2021;48(9):844-850.

31. Jones P, Binns D, Chang HY, et al. InterProScan 5: genome-scale protein function classification. *Bioin*formatics.2014;30(9):1236-1240.

32. Blum M, Chang HY, Chuguransky S, et al. The InterPro protein families and domains database: 20 years on. *Nucleic Acids Res*.2021;49(D1):D344-d354.

33. Manriquez-Sandoval E, Fried SD. DomainMapper: Accurate domain structure annotation including those with non-contiguous topologies. *Protein Sci.* 2022;31(11):e4465.

34. Hallgren J, Tsirigos KD, Pedersen MD, et al. DeepTMHMM predicts alpha and beta transmembrane proteins using deep neural networks. *bioRxiv*. 2022:2022.2004.2008.487609.

35. Bitard-Feildel T, Lamiable A, Mornon JP, Callebaut I. Order in Disorder as Observed by the "Hydrophobic Cluster Analysis" of Protein Sequences. *Proteomics.* 2018;18(21-22):e1800054.

36. Callebaut I, Labesse G, Durand P, et al. Deciphering protein sequence information through hydrophobic cluster analysis (HCA): current status and perspectives. *Cell Mol Life Sci.* 1997;53(8):621-645.

37. Bruley A, Bitard-Feildel T, Callebaut I, Duprat E. A sequence-based foldability score combined with AlphaFold2 predictions to disentangle the protein order/disorder continuum. *Proteins*.2023;91(4):466-484.

38. Bruley A, Mornon JP, Duprat E, Callebaut I. Digging into the 3D Structure Predictions of AlphaFold2 with Low Confidence: Disorder and Beyond. *Biomolecules*. 2022;12(10).

39. Finn RD, Clements J, Eddy SR. HMMER web server: interactive sequence similarity searching. *Nucleic Acids Res.* 2011;39(Web Server issue):W29-37.

40. Larralde M, Zeller G. PyHMMER: a Python library binding to HMMER for efficient sequence analysis. *Bioinformatics*. 2023;39(5).

41. Steinegger M, Söding J. MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nat Biotechnol*.2017;35(11):1026-1028.

42. Gibbons TR, Mount SM, Cooper ED, Delwiche CF. Evaluation of BLAST-based edge-weighting metrics used for homology inference with the Markov Clustering algorithm. *BMC Bioinformatics*. 2015;16:218.

43. Hagberg A, Swart P, S Chult D. Exploring network structure, dynamics, and function using networkx. Conference: SCIPY 08; August 21, 2008; Pasadena; 2008; United States.

44. Blondel VD, Guillaume J-L, Lambiotte R, Lefebvre E. Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment.* 2008;2008(10):P10008.

45. Shannon P, Markiel A, Ozier O, et al. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res.* 2003;13(11):2498-2504.

46. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol.* 2013;30(4):772-780.

47. Waterhouse AM, Procter JB, Martin DM, Clamp M, Barton GJ. Jalview Version 2–a multiple sequence alignment editor and analysis workbench. *Bioinformatics*. 2009;25(9):1189-1191.

48. Crooks GE, Hon G, Chandonia JM, Brenner SE. WebLogo: a sequence logo generator. *Genome Res.* 2004;14(6):1188-1190.

49. Pettersen EF, Goddard TD, Huang CC, et al. UCSF ChimeraX: Structure visualization for researchers, educators, and developers. *Protein Sci.* 2021;30(1):70-82.

50. Tunyasuvunakool K, Adler J, Wu Z, et al. Highly accurate protein structure prediction for the human proteome. *Nature*.2021;596(7873):590-596.

51. van Kempen M, Kim SS, Tumescheit C, et al. Fast and accurate protein structure search with Foldseek. *Nat Biotechnol.* 2023.

52. Robert X, Gouet P. Deciphering key features in protein structures with the new ENDscript server. *Nucleic Acids Res.* 2014;42(Web Server issue):W320-324.

53. Andersson M, Mattle D, Sitsel O, et al. Copper-transporting P-type ATPases use a unique ion-release pathway. *Nat Struct Mol Biol*.2014;21(1):43-48.

54. Gourdon P, Liu XY, Skjørringe T, et al. Crystal structure of a copper-transporting PIB-type ATPase. *Nature*.2011;475(7354):59-64.

55. Toyoshima C, Nakasako M, Nomura H, Ogawa H. Crystal structure of the calcium pump of sarcoplasmic reticulum at 2.6 A resolution. *Nature*. 2000;405(6787):647-655.

56. Smith AT, Smith KP, Rosenzweig AC. Diversity of the metal-transporting P1B-type ATPases. J Biol Inorg Chem. 2014;19(6):947-960.

57. Thomas C, Aller SG, Beis K, et al. Structural and functional diversity calls for a new classification of ABC transporters. *FEBS Lett.* 2020;594(23):3767-3775.

58. Leskelä S, Kontinen VP, Sarvas M. Molecular analysis of an operon in Bacillus subtilis encoding a novel ABC transporter with a role in exoprotein production, sporulation and competence. *Microbiology (Reading)*. 1996;142 (Pt 1):71-77.

59. Sigal YJ, McDermott MI, Morris AJ. Integral membrane lipid phosphatases/phosphotransferases: common structure and diverse functions. *Biochem J.* 2005;387(Pt 2):281-293.

60. Stukey J, Carman GM. Identification of a novel phosphatase sequence motif. *Protein Sci.* 1997;6(2):469-472.

61. Dillon DA, Wu WI, Riedel B, Wissing JB, Dowhan W, Carman GM. The Escherichia coli pgpB gene encodes for a diacylglycerol pyrophosphate phosphatase activity. *J Biol Chem.* 1996;271(48):30548-30553.

62. Zhao J, An J, Hwang D, et al. The Lipid A 1-Phosphatase, LpxE, Functionally Connects Multiple Layers of Bacterial Envelope Biogenesis. *mBio.* 2019;10(3).

63. Miller DJ, Jerga A, Rock CO, White SW. Analysis of the Staphylococcus aureus DgkB structure reveals a common catalytic mechanism for the soluble diacylglycerol kinases. *Structure*.2008;16(7):1036-1046.

64. Bakali MA, Nordlund P, Hallberg BM. Expression, purification, crystallization and preliminary diffraction studies of the mammalian DAG kinase homologue YegS from Escherichia coli. Acta Crystallogr Sect F Struct Biol Cryst Commun. 2006;62(Pt 3):295-297.

65. Borgstahl GE, Parge HE, Hickey MJ, Beyer WF, Jr., Hallewell RA, Tainer JA. The structure of human mitochondrial manganese superoxide dismutase reveals a novel tetrameric interface of two 4-helix bundles. *Cell*. 1992;71(1):107-118.

66. Kleiger G, Grothe R, Mallick P, Eisenberg D. GXXXG and AXXXA: common alpha-helical interaction motifs in proteins, particularly in extremophiles. *Biochemistry*. 2002;41(19):5990-5997.

67. Beamer LJ, Carroll SF, Eisenberg D. Crystal structure of human BPI and two bound phospholipids at 2.4 angstrom resolution. *Science*.1997;276(5320):1861-1864.

68. Kleiger G, Beamer LJ, Grothe R, Mallick P, Eisenberg D. The 1.7 A crystal structure of BPI: a study of how two dissimilar amino acid sequences can adopt the same fold. *J Mol Biol*.2000;299(4):1019-1034.

69. Kopec KO, Alva V, Lupas AN. Bioinformatics of the TULIP domain superfamily. *Biochem Soc Trans.* 2011;39(4):1033-1038.

70. Koehler Leman J, Szczerbiak P, Renfrew PD, et al. Sequence-structure-function relationships in the microbial protein universe. *Nat Commun.* 2023;14(1):2351.

71. Jordan IK, Natale DA, Koonin EV, Galperin MY. Independent evolution of heavy metal-associated domains in copper chaperones and copper-transporting atpases. J Mol Evol. 2001;53(6):622-633.

72. Wang K, Sitsel O, Meloni G, et al. Structure and mechanism of Zn2+-transporting P-type ATPases. *Nature*. 2014;514(7523):518-522.

73. Mattle D, Sitsel O, Autzen HE, Meloni G, Gourdon P, Nissen P. On allosteric modulation of P-type Cu(+)-ATPases. J Mol Biol.2013;425(13):2299-2308.

74. Biemans-Oldehinkel E, Doeven MK, Poolman B. ABC transporter architecture and regulatory roles of accessory domains. *FEBS Lett.* 2006;580(4):1023-1035.

75. Sikkema HR, van den Noort M, Rheinberger J, et al. Gating by ionic strength and safety check by cyclic-di-AMP in the ABC transporter OpuA. *Sci Adv.* 2020;6(47).

76. Hwang TC, Braakman I, van der Sluijs P, Callebaut I. Structure basis of CFTR folding, function and pharmacology. J Cyst Fibros.2023;22 Suppl 1:S5-s11.

77. Coleman JA, Quazi F, Molday RS. Mammalian P4-ATPases and ABC transporters and their role in phospholipid transport. *Biochim Biophys Acta*. 2013;1831(3):555-574.

78. Lyons JA, Timcenko M, Dieudonné T, Lenoir G, Nissen P. P4-ATPases: how an old dog learnt new tricks - structure and mechanism of lipid flippases. *Curr Opin Struct Biol.* 2020;63:65-73.

79. Wong LH, Gatta AT, Levine TP. Lipid transfer proteins: the lipid commute via shuttles, bridges and tubes. *Nat Rev Mol Cell Biol*.2019;20(2):85-101.

80. Wong LH, Levine TP. Tubular lipid binding proteins (TULIPs) growing everywhere. *Biochim Biophys Acta Mol Cell Res*.2017;1864(9):1439-1449.

81. Corbalan-Garcia S, Gómez-Fernández JC. Signaling through C2 domains: more than one lipid target. *Biochim Biophys Acta*.2014;1838(6):1536-1547.

82. Lemmon MA. Membrane recognition by phospholipid-binding domains. *Nat Rev Mol Cell Biol.* 2008;9(2):99-111.

83. Malcova I, Bumba L, Uljanic F, Kuzmenko D, Nedomova J, Kamanova J. Lipid binding by the N-terminal motif mediates plasma membrane localization of Bordetella effector protein BteA. J Biol Chem.2021;296:100607.

84. Varela-Chavez C, Blondel A, Popoff MR. Bacterial intracellularly active toxins: Membrane localisation of the active domain. *Cell Microbiol.* 2020;22(7):e13213.

85. López-Lara IM, Geiger O. Bacterial lipid diversity. *Biochim Biophys Acta Mol Cell Biol Lipids*. 2017;1862(11):1287-1299.

86. Petroutsos D, Amiar S, Abida H, et al. Evolution of galactoglycerolipid biosynthetic pathways-from cyanobacteria to primary plastids and from primary to secondary plastids. *Prog Lipid Res.* 2014;54:68-85.

87. Wada H, Murata N. Membrane lipids in cyanobacteria. In: Siegenthaler P-A, Murata N, eds. *Lipids in photosynthesis*. Dordrecht, The Netherlands: Kluwer Academic Publishers; 1998:65-81.

88. Boudière L, Michaud M, Petroutsos D, et al. Glycerolipids in photosynthesis: composition, synthesis and trafficking. *Biochim Biophys Acta*. 2014;1837(4):470-480.

89. Wada H, Murata N. The essential role of phosphatidylglycerol in photosynthesis. *Photosynth Res.* 2007;92(2):205-215.

90. Kóbori TO, Uzumaki T, Kis M, et al. Phosphatidylglycerol is implicated in divisome formation and metabolic processes of cyanobacteria. J Plant Physiol. 2018;223:96-104.

91. Terrapon N, Weiner J, Grath S, Moore AD, Bornberg-Bauer E. Rapid similarity search of proteins using alignments of domain arrangements. *Bioinformatics*. 2014;30(2):274-281.

92. Buchan DWA, Jones DT. Learning a functional grammar of protein domains using natural language word embedding techniques. *Proteins*. 2020;88(4):616-624.