

A catalog of the diversity and ubiquity of bacterial microcompartments

Markus Sutter ^{1,2}, Matthew R. Melnicki², Frederik Schulz ³, Tanja Woyke ³ & Cheryl A. Kerfeld ^{1,2,4}✉

Bacterial microcompartments (BMCs) are organelles that segregate segments of metabolic pathways which are incompatible with surrounding metabolism. BMCs consist of a selectively permeable shell, composed of three types of structurally conserved proteins, together with sequestered enzymes that vary among functionally distinct BMCs. Genes encoding shell proteins are typically clustered with those for the encapsulated enzymes. Here, we report that the number of identifiable BMC loci has increased twenty-fold since the last comprehensive census of 2014, and the number of distinct BMC types has doubled. The new BMC types expand the range of compartmentalized catalysis and suggest that there is more BMC biochemistry yet to be discovered. Our comprehensive catalog of BMCs provides a framework for their identification, correlation with bacterial niche adaptation, experimental characterization, and development of BMC-based nanoarchitectures for biomedical and bioengineering applications.

¹Environmental Genomics and Systems Biology and Molecular Biophysics and Integrative Bioimaging Divisions, Lawrence Berkeley National Laboratory, Berkeley, CA, USA. ²MSU-DOE Plant Research Laboratory, Michigan State University, East Lansing, MI, USA. ³DOE Joint Genome Institute, Lawrence Berkeley National Laboratory, Berkeley, CA, USA. ⁴Department of Biochemistry and Molecular Biology, Michigan State University, East Lansing, MI, USA. ✉email: ckerfeld@lbl.gov

Bacterial microcompartments (BMCs) are metabolic organelles that consist entirely of protein; a modular shell surrounds an enzymatically active core, with the shell functioning as a semipermeable membrane for substrates and products (Fig. 1a). The first type of BMC discovered were carboxysomes, they were observed in electron micrographs as polyhedral structures in Cyanobacteria¹. Carboxysomes enhance CO₂ fixation² in all cyanobacteria and some chemoautotrophs by encapsulating RuBisCO together with carbonic anhydrase to concentrate the substrate CO₂ (Fig. 1b). Much later, similar structures were observed in heterotrophs, however only when grown in the presence of the substrate of those BMCs, ethanolamine or 1,2-propanediol³. DNA sequencing confirmed that the shell proteins of those BMCs are similar to those of the carboxysomes, a fact that has enabled finding a multitude of BMCs with the advent of genomic sequencing^{4,5}. The majority of BMCs are catabolic and are known as metabolosomes^{4,6}. Many diverse types share a common core chemistry of a signature enzyme^{4,7}, that generates an aldehyde, an ubiquitous pfam00171 aldehyde dehydrogenase (AldDh) to oxidize it as well as a phosphotransacylase (PTAC)⁸ to generate an acyl-phosphate and an alcohol dehydrogenase (AlcDh) for cofactor regeneration (Fig. 1c)⁷. BMCs protect the cytosol from toxic intermediates and enhance catalysis by co-localizing enzymes and concentrating substrates⁹. Targeting of enzymes into the lumen of many BMCs,

including the beta-carboxysome¹⁰, typically proceeds via encapsulation peptides (EPs), which consist of a ~15–20 amino acid amphipathic alpha-helix that is connected to the N- or C-termini of cargo proteins via a flexible linker¹¹.

BMCs have a polyhedral, often icosahedral shape, and typically range from about 40–200 nm in diameter⁹. Structures of model shells confirmed icosahedral symmetry with pentagons occupying 12 vertices and hexagons forming the facets^{12–15} (Fig. 1a). The pentagons (BMC-P) are formed by five subunits of the pfam03319 fold and have the shape of a truncated pyramid¹⁶ (Supplementary Fig. 1a). BMC-H proteins contain the pfam00936 domain and form almost perfect hexagons with a diameter of about 65 Å and with a concave and convex side¹⁷ (Supplementary Fig. 1b). A central pore allows for passage of substrates and products¹⁷. The facets of shells also contain BMC-T proteins, which are genetic fusions of two pfam00936 domains that form trimers that in size and shape resemble the hexamers. Variations of the pfam00936 domain that are the result of a circular permutation of the primary structure expand the number of distinct hexagonal building blocks to five (Fig. 1a, Supplementary Fig. 1b–f).

In 2014, a comprehensive bioinformatic survey identified 357 BMC loci representative of 30 types/subtypes across 23 bacterial phyla⁴. This cataloging of BMC loci within genomic databases was feasible, because both pfam03319 and pfam00936 protein

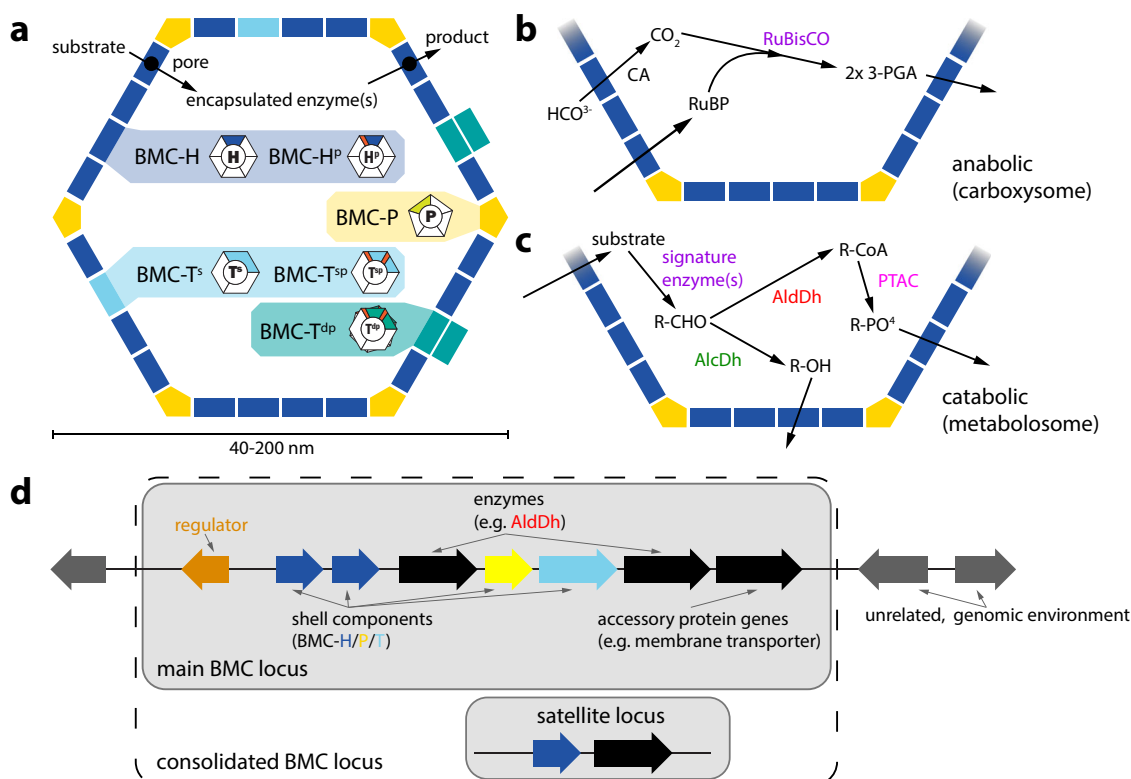


Fig. 1 Generalized BMC structure, function, and chromosomal organization of component genes. **a** Overview of a BMC shell and types of shell protein components. BMC-P: pentamer, pfam03319 domain (Supplementary Fig. 1a); BMC-H: hexamer, pfam00936 domain (Supplementary Fig. 1b); BMC-HP: circularly permuted variant of BMC-H with two secondary structure elements translocated from the C- to the N-terminus (Supplementary Fig. 1e); BMC-T^S: standard trimer, fusion of two pfam00936 domains (Supplementary Fig. 1c); BMC-T^{SP}: permuted BMC-T^S, each pfam00936 domain contains a circular permutation as in BMC-HP (Supplementary Fig. 1f); BMC-T^{DP}: a permuted BMC-T variant in which two trimers dimerize across their concave faces to form an interior chamber (Supplementary Fig. 1d). The pores in BMC-T^{DP} trimers are relatively large (14 Å in diameter), gated by conformational changes in the surrounding sidechains and are predicted to serve as conduits for larger metabolites^{65–67}. **b, c** Simplified reactions of anabolic (**b**) and catabolic (**c**) BMCs. CA carbonic anhydrase, RuBP ribulose 1,5-bisphosphate, 3-PGA: 3-phosphoglycerate, AlcDh alcohol dehydrogenase, AldDh aldehyde dehydrogenase, PTAC phosphotransacylase, CoA coenzyme A. **d** Typical BMC locus consisting of genes for shell proteins (blue, cyan, yellow), enzymes (black, colored according to enzyme type in detailed diagrams of Supplementary Data 2 and Supplementary Data 3), regulators (orange) and ancillary proteins such as cell membrane transporters for substrates; The combination of main and satellite BMC locus is termed a consolidated BMC locus.

fold are unique to BMCs and are often encoded together with their cargo proteins in chromosomal loci. Within the last six years, the number of metagenomic datasets has increased nearly an order of magnitude¹⁸. In addition to the microbiome data that have become available, the extraction of tens of thousands of genomes from the uncultivated majority of microorganisms has been facilitated through recent advances in genome-resolved metagenomics¹⁹.

In this work, we have mined these massive new datasets and compiled a database of more than 7000 BMC loci that cluster into 68 BMC types or subtypes, including 29 new functional BMC types or subtypes. BMC loci are widespread, now evident in 45 phyla across the bacterial tree of life. Collectively our results show that the known BMC functional diversity and distribution at the phylum level has essentially doubled in the last six years and foregrounds the widespread occurrence of bacterial organelles.

Results

BMC shell protein and locus analysis. We hypothesized that BMC diversity has greatly expanded since the previous survey (2014) due to growth in genome sequencing of microbial diversity and in particular that of uncultivated clades. We compiled and curated an in-house dataset of all putative BMC loci based on the UniProt Knowledgebase (UniProtKB) with data as recent as March 2020. After retrieving all available BMC shell protein sequences, we collected the sequences for neighboring genes encoded within 12 ORFs from any shell protein gene, which covers all BMC locus-related genes in previously known loci and, in retrospect, in all new BMC loci. Contiguous genes were classified as a “main locus” when they contained at least one BMC-H and one BMC-P gene, to distinguish it from “satellite loci”⁴. Satellite loci were combined with the main locus to form a consolidated locus (Fig. 1d).

The previous comprehensive survey⁴ relied on pfam co-occurrence. Many of the new BMC loci would be difficult to categorize by this method, because they contained few recognizably type-specific proteins other than the shell components (unlike the well-characterized PDU and EUT loci with around 20 total gene products; Fig. 2a). We therefore sought to improve BMC type clustering by subclassifying all pfam00936 and pfam03319 shell proteins using a phylogenomic approach²⁰. Trees were built for representative sequences from each of the six shell protein types: BMC-H, BMC-P, BMC-T^s, BMC-T^{dp}, BMC-H^p, and BMC-T^{sp} (Figs. 1a, 3). Subclades for each shell type were identified visually, usually containing a long internal stem, and were each assigned a unique color name chosen from the xkcd color survey (<https://xkcd.com/color/rgb/>), using names from related color families for adjacent subclades within each major clade (for high-resolution trees with full annotations, see Supplementary Data 1). The sequences comprising each of the colored subclades were then used to calculate a profile Hidden Markov Model (HMM)²¹ for each color group and combined with HMMs derived from proteins common to BMC loci that were clustered by protein pfam information (Supplementary Fig. 2a). Decoupling the identity of a shell protein from a specific BMC type allowed us to make unbiased observations of shell proteins that are similar, despite being constituents of functionally distinct BMCs. We find that for many loci, the component shell proteins are drawn from across the tree, revealing functional and evolutionary relationships among distinct BMC types (described below).

The combination of detailed shell protein, enzyme, and accessory protein HMMs allowed us to cluster loci into distinct BMC types and subtypes (Supplementary Fig. 2a, see “Methods” for details). Subtypes generally share most enzymatic components

but differ in gene order and/or shell protein content. For example, the sugar-phosphate utilization (SPU loci) can be separated into seven distinct subtypes (Supplementary Fig. 2b). For the naming of the loci we adopted and expanded the nomenclature from Axen et al.⁴ (Fig. 2c). New BMCs were named either by a distinguishing feature such as predominant occurrence in a certain taxon of organisms (e.g., ACI for Acidobacteria) or the organism that gave rise to a model system (e.g., HO for *Haliangium ochraceum*), or putative class of BMC substrate (e.g., ARO for aromatic substrate). In absence of any potentially defining feature, BMCs were classified under broad groupings such as Metabolosomes of unknown function (MUF), Metabolosomes with an incomplete core (MIC) or, for BMCs that lack an AldDh, BMC of unknown function (BUF) (Fig. 2c). In this analysis, we have added 29 new major BMC types or subtypes as well as 10 new subtypes for several established loci (Fig. 2a). The total number of loci is more than 7000, an increase of almost 20-fold from the 357 loci listed in the study from 2014⁴. Representative locus diagrams and a short description of the common components for a total of 68 types are listed in Supplementary Data 2.

Overview of BMC distribution and characteristics. BMCs occur widely across the bacterial domain (Fig. 4). Compared to 23 phyla with BMCs in 2014⁴, we now find BMCs in 45 out of 83 phyla and proteobacterial classes with genomic sequence data in IMG¹⁸ (counting Patescibacteria as a single phylum). Some BMCs are confined to a single phylum (e.g., ACI, ARO, BUF2) but most are found across multiple phyla (Fig. 4). Our analysis revealed a large BMC functional diversity in certain phyla, such as Proteobacteria, Actinobacteria, and Firmicutes (Fig. 4), reflecting the importance of metabolic flexibility across disparate niches. To get an estimate of the overall prevalence of BMCs we performed a shell protein HMM search against IMG/M¹⁸ and GEM¹⁹ genomes and found hits in 20% of them.

A core enzyme of many functionally distinct metabolosomes is the AldDh (Fig. 2a). A phylogenetic tree of representative sequences reveals that the AldDh is specific to BMC functional type (Fig. 5). AldDh from BMCs with similar substrates cluster on the tree and, accordingly, can be used in the prediction of potential substrates for unknown BMCs by looking at the closest AldDh homologs. AldDh typically have an EP on either the N- and C- termini; strikingly the two major branches of the tree also are distinct in the location of the EP extensions (Fig. 5); type I have EP at their C-terminus, type II at the N-terminus. The most parsimonious interpretation of these data is that the acquisition of a sequence extension that serves to facilitate encapsulation is an ancient innovation that arose independently twice in the evolution of BMCs.

In addition to enzymes, BMC loci also contain genes for functions that support the expression and activity of the organelle like transcriptional regulators and cell membrane transporters for the substrate. Like AldDh, these gene products also can provide clues as to the function of a BMC, by comparison to non-BMC-related homologs. For example, for one of the new BMC types that we identified here involves an aromatic substrate (ARO), a SWISS-MODEL search²² with the protein sequence of the regulator reveals the top characterized hit as a regulator of catechol degradation, which is consistent with the enzymes found in that BMC locus (Supplementary Data 2).

When comparing the shell protein inventories across the BMC types (Fig. 2a), both unexpected differences and new patterns emerge. On average, BMC loci contain 1.7 BMC-P, 2.2 BMC-H, 0.4 BMC-Hp, 0.5 BMC-T^s, 0.4 BMC-T^{sp} and 0.3 BMC-T^{dp}. The distribution reveals that BMC-P most commonly occurs singly or

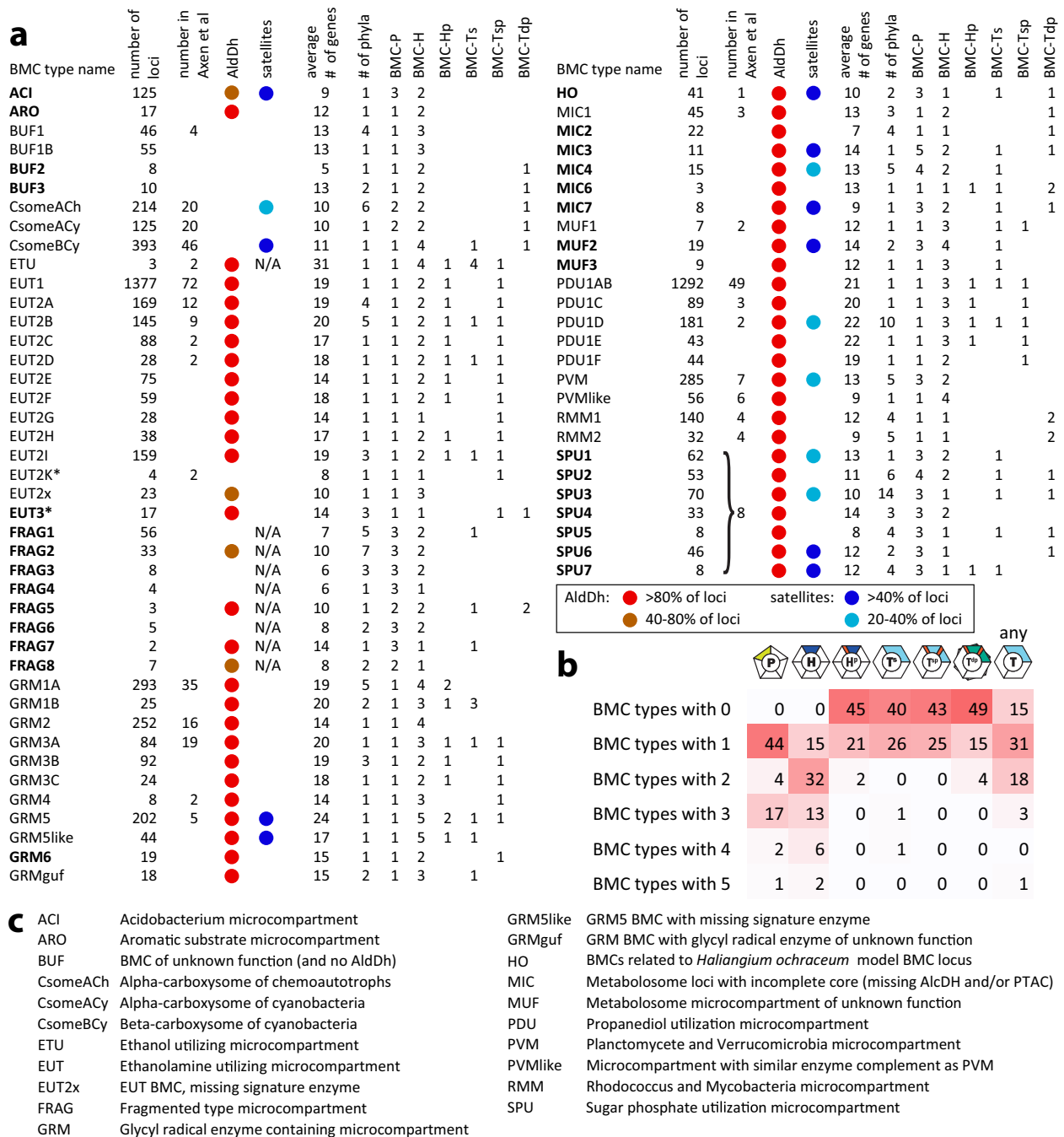


Fig. 2 Overview of BMC types and shell protein content found across all bacteria. **a** Table of number of loci for each BMC type in this study and in Axen et al.⁴, AldDh occurrence, the prevalence of satellite loci, the average number of genes in the locus, number of observed phyla, and number of each type of shell protein. Major new BMC types or subtypes identified in this study in bold. The asterisk denotes the name change of EUT3 (Axen et al.⁴) to EUT2K because the present analysis indicates EUT2K is not a major new type. **b** Numbers of each type of shell protein across BMC types. **c** Explanation of BMC type abbreviations.

as three paralogs (Fig. 2b). Correlating the BMC-P found in the loci that have three copies with their location on the BMC-P tree (Fig. 3a), we find a “BMC-P triplet” pattern: one member each from the gray and orange major clades and a third member from one of the other clades (Supplementary Fig. 3). Many BMC loci also contain multiples of BMC-H and for a specific BMC type, those do not necessarily cluster together on a phylogenetic tree. A reason for this could be that the different paralogs fulfill a function that is shared across BMC types. The wide variety of different compositions of the BMC shell highlights its modular

construction from building blocks that have the same size and shape, but different permselectivities.

New and expanded variants of BMC loci. One of the most prominent expanded BMC types is the SPU (sugar phosphate utilization), first noted in 2014 through the identification of only eight representatives⁴. We now find more than 280 members distributed across 26 bacterial phyla (Fig. 4), establishing SPU as one of the most prevalent BMC types. Our clustering has

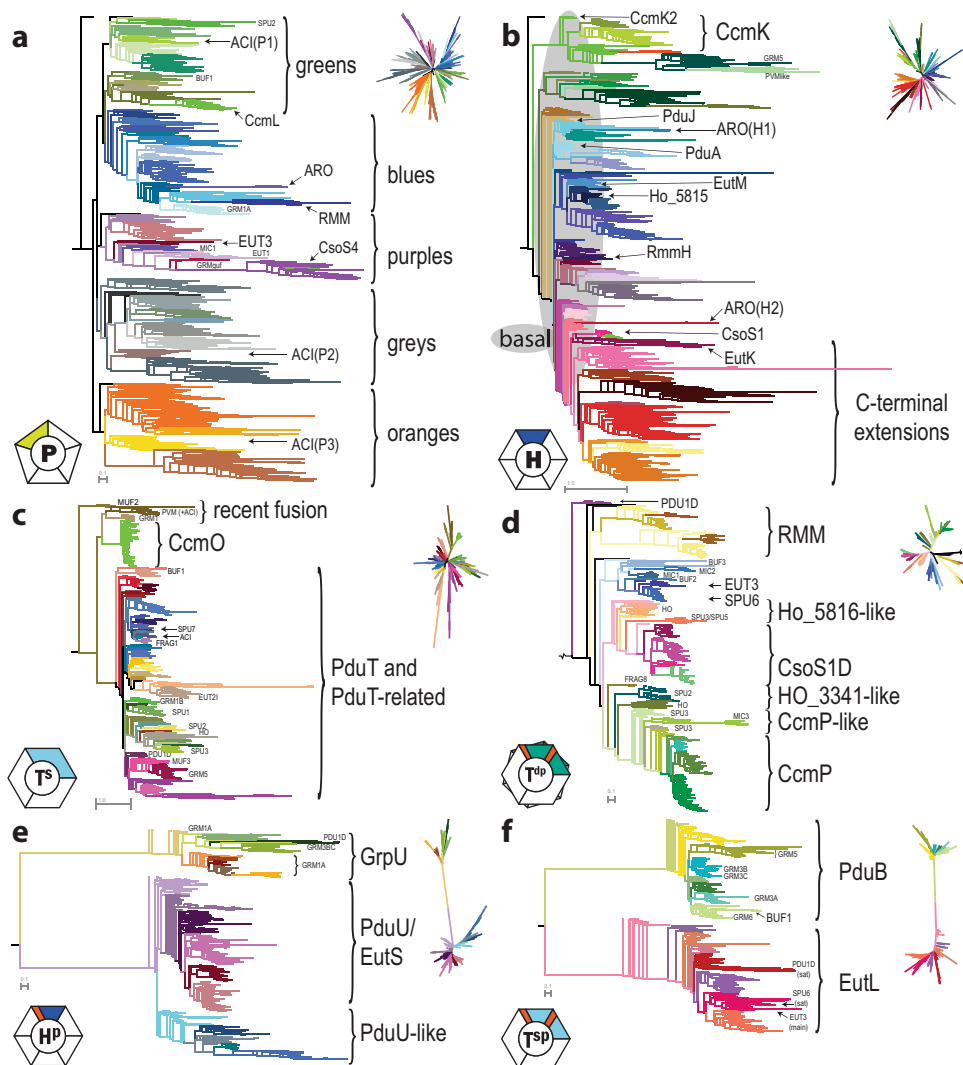


Fig. 3 Phylogenetic maximum likelihood trees for the six types of shell protein. a BMC-P, **b** BMC-H with the region of basal hexamers encompassed by gray shading, **c** BMC-T^s, **d** BMC-T^{dp}, **e** BMC-HP and **f** BMC-T^{sp}. Representative sequences were selected by removing redundancy. Full hi-resolution images of each tree with annotation labels for all terminal nodes are available as Supplementary Data 1. All caps labels refer to the predominantly associated locus with a certain BMC type, mixed case labels denote a specific, previously characterized protein found in that clade. For BMC-P the major clades are labeled by color. Unrooted versions of the trees are shown in the top right corner of each panel.

identified seven distinct subtypes (Supplementary Fig. 2b, Supplementary Data 2) that all share two sugar phosphate processing enzymes (pfam01791, pfam02502). The DeoC-type pfam01791 aldolase converts 2-deoxy-D-ribose to glyceraldehyde-3-phosphate and acetaldehyde which can then be processed by the AldDh to acetyl-CoA. The AldDh of SPU6 is close to EUT3 and the other SPU types are on the same major branch as EUT1, which both process acetaldehyde (Fig. 5). The pfam02502 is of the RpiB type²³ that isomerizes ribose-5-phosphate. Collectively, these data suggest potential function of this BMC type is metabolizing the products of DNA degradation, presumably from the detritus ubiquitously available in diverse environments.

The HO BMC from *Haliangium ochraceum* has been the primary model system for structural studies of the shell^{12,14,24}. We have identified similar loci in 40 other genomes but the function of this organelle remains enigmatic; the HO AldDh is most closely related to those of SPU5 and SPU7 (Fig. 5), all loci encode the characteristic BMC-P triplet, and share similar types of BMC-T^{dp} (Fig. 3d). Some genomes containing the HO loci have a pfam01791 sugar processing enzyme in a different

genomic location and some of these orthologs contain an EP-like N-terminal extension and could therefore be encapsulated in aHO BMC. A similar function as SPU BMCs, the catabolism of nucleic acid, is consistent with their presence in Myxobacteria that are known for degrading biomass derived from dead cells²⁵.

Another functional type that has substantially increased in membership are the PVMs²⁶ now with 285 representatives (Fig. 2), as compared to seven found in 2014⁴. This reflects the increased attention to Verrucomicrobia species for their role in the global carbon cycle as degraders of complex algal and bacterial cell wall polysaccharides^{27,28}. For example, *Lentimonas* species devote 4% of their proteome to the degradation of fucoidan, the major cell wall polysaccharide of brown algae into fucose, which is catabolized in the PVM BMC²⁸. The PVM-like BMC that shares the PVM aldolase signature enzyme that processes sugar derivatives also gained 50 new members and can be found in several bacteria of the human gut microbiome²⁹. We expect more PVM and PVM-like BMCs to be discovered with increased attention on sequencing bacteria that degrade complex

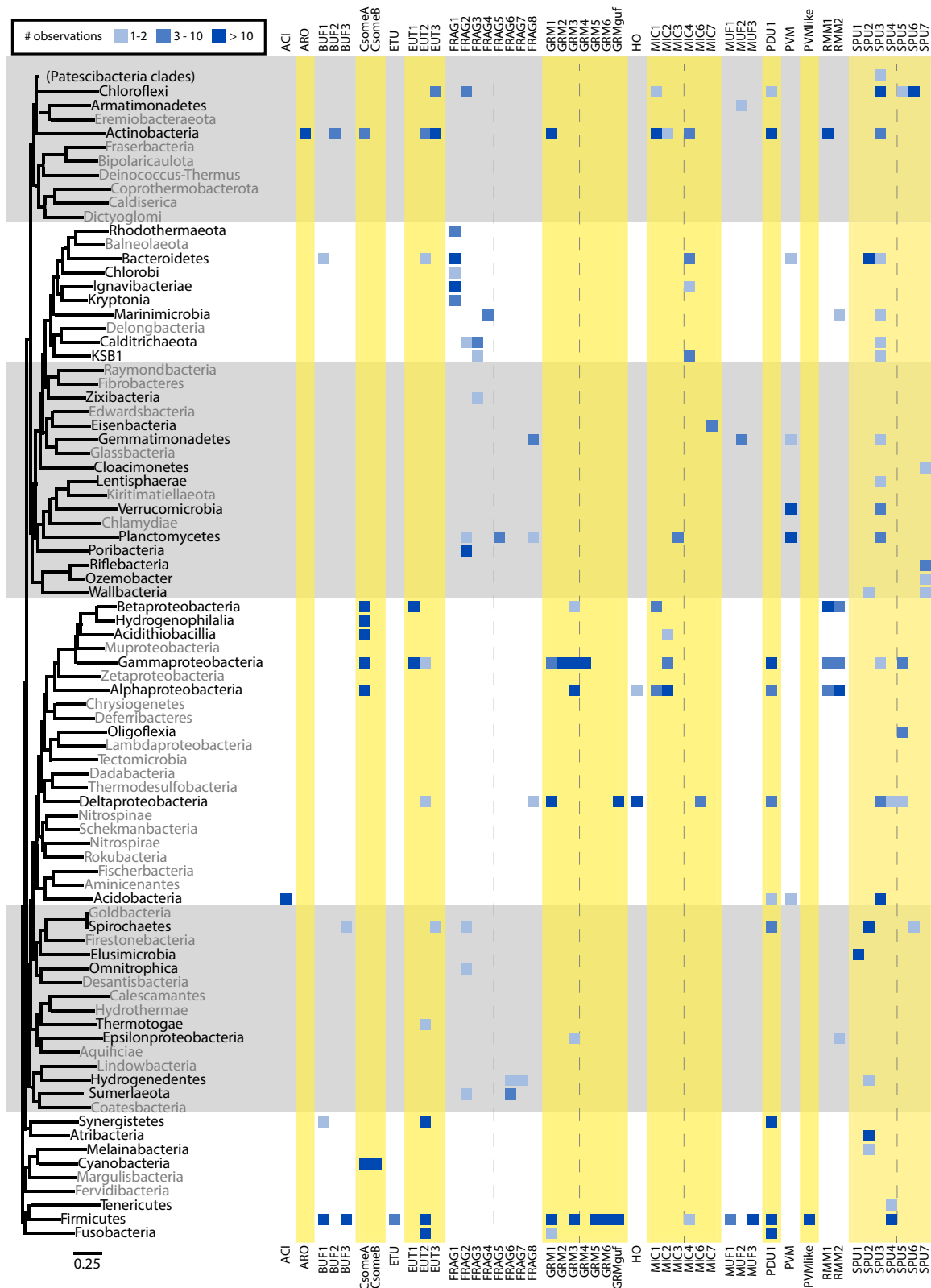


Fig. 4 Distribution of BMC types in 45 bacterial phyla. Phylogenetic tree representing all phylum-level taxonomic groups was generated by aligning a set of 56 different marker proteins. Phyla lacking BMC loci shown in gray and different numbers of BMCs in blue. Alternate yellow/gray blocks and dashed lines are added for visual guidance.

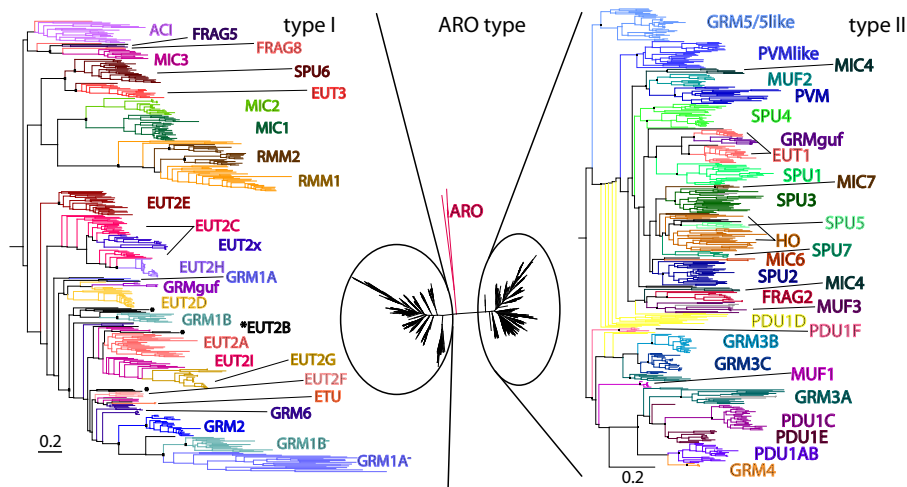


Fig. 5 Phylogenetic tree of pfam00171 aldehyde dehydrogenases from BMC loci. Sequence redundancy was reduced by limiting to 30 sequences per BMC type (based on percent identity). Three major groups were identified: type I, which is dominated by EUT2, belongs to IPR013357, and contains a C-terminal EP (left); the ARO type outlier clade; and type II, which belongs to IPR012408, contains an N-terminal EP, and contains a larger diversity of BMC types, relative to type I AldDs. Branch lengths are scaled by the number of substitutions per site. Bootstrap values for important nodes are represented as black squares (above 50%). GRM1A and GRM1B loci contain a putatively inactive second copy (marked with a “-” superscript). EUT2B AldDh marked with an asterisk are found in different basal locations.

polysaccharides, and we find their shell protein components prevalent in searches of metagenomes from environmental samples (see below).

Our analysis discovered several completely new BMC types. One, we named ARO, for its predicted aromatic substrate, is found in the Micromonosporales and Pseudonocardiales orders of Actinobacteria. The ARO locus contains two pfam02900 ring-opening oxygenases and a set of enzymes related to the degradation of aromatic aldehyde compounds (Supplementary Data 2). A possible initial substrate is 2-aminophenol based on the assignment of related pfam00171 AldDh as aminomuconate-semialdehyde dehydrogenases. Unlike most other AldDs, the ARO AldDh does not contain a detectable EP, and constitutes its own group falling between both major groups on the phylogenetic tree (Fig. 5). Likewise, the ARO shell protein composition is among the simplest observed: consisting of one BMC-P and two distinct types of BMC-H. The three ARO shell proteins are all found in late-branching subclades that are strongly divergent from other shell proteins (Fig. 3a, b), suggesting that the predicted catabolism of aromatic compounds and the involvement of an AldDh remote from others found in BMCs has imposed a distinctive function on these shell proteins that is under evolutionary constraint.

Another major new type with 125 members, ACI, is found exclusively in Acidobacteria, primarily in only two of the 26 Subdivisions, Subdivision 4 and Subdivision 6. While Acidobacteria are found across differing ecosystems³⁰ they are of particular relevance to the soil environment, as they can comprise up to 60% of the soil bacterial community³¹. The locus contains a pfam00596 class II aldolase with a C-terminal EP, a hydroxyacid dehydrogenase (pfam00389/02826), a triplet of BMC-P proteins, at least one BMC-H and several proteins of unknown function (Supplementary Data 2). Its AldDh is consistently found on a satellite locus and is phylogenetically similar only to MIC3, another uncharacterized BMC type (Fig. 5). The pfam00596 aldolase is also observed in the PVM and GRM5 types that process L-fucose and L-rhamnulose phosphate so a function related to carbohydrate degradation can be proposed.

We discovered several enigmatic new BMC types that lack a defined locus organization but are able to be grouped based solely

on shell protein composition. These FRAG BMCs are composed of shell protein genes encoded by as many as six different genomic locations. One common feature of FRAG BMCs is the presence of a BMC-P triplet (Fig. 3a, Supplementary Fig. 3). Most of these genes are remote from genes for known enzymes, although a few have proximal AldDs (Fig. 5, Supplementary Data 2), that map to uncharacterized AldDs in different parts of the tree. FRAG BMCs are found in diverse organisms, including the Ignavibacteriae and closely related phyla, the Gemmatimonadetes, the Planctomycetes, and a large number of candidate phyla (Fig. 4). In addition, we have found several more distinct BMC types of unknown function: BUF2, BUF3, GRM6, MIC2-7, MUF2, and MUF3 (Supplementary Data 2). There is a large number of new metabolosome loci that encode a pfam00171 AldDh yet do not have an obvious signature enzyme that generates that aldehyde (ACI, FRAG1-7, HO, MIC2-7, MUF2-3; Fig. 2a, c). Assignment of function for these and other functionally cryptic loci could be approached by combining genetic, biochemical studies and metabolic arrays focused on the AldDH and the shell proteins, as done for the identification of the fucose/rhamnose substrate for PVM BMCs²⁶.

BMC type co-occurrence and horizontal gene transfer. About 80% of BMC-containing genomes encode only one locus, but a substantial number encode two (20%) or more (2%) loci. The most frequent co-occurring pair are the PDU and EUT loci (Supplementary Table 1), providing catabolism for ethanolamine and 1,2-propanediol. In *Salmonella enterica*, this combination has been shown to be regulated mutually exclusively to prevent the formation of mixed BMCs³². Most of the genomes encoding three or more BMC types are combinations of EUT, PDU, and GRM (Supplementary Table 2). The most extreme examples of the potential to form multiple metabolic organelles is the genome of the Firmicute *Maledivibacter halophilus* that contains six BMC loci: two EUT2B and one each of PDU1D, GRM1A, BUF1 and BUF3. This organism is found in anoxic hypersaline sediments³³. All of the organisms with three or more BMC types are either Proteobacteria or Firmicutes, predominantly from the Enterobacterales or Clostridiales orders. Members of these orders

are various human and other animal pathogens, and are likewise abundant in aquatic and soil ecosystems. The ability to form multiple functionally distinct BMCs confers metabolic potential and flexibility; for example, the prevalence of multiple BMCs in the Firmicutes (*Clostridiaceae*, *Ruminococcaceae*) likely contributes to their ability to break down of complex polysaccharides³⁴.

BMC loci are genetic modules, a compact organization of the structural, regulatory, and ancillary components necessary for BMC function. As such, they are an ideally suited for horizontal gene transfer (HGT). One example of a possibly recent HGT event involves a EUT2I locus in *Oceanotoga teriensis*, the only BMC instance in the Thermotogae phylum. About 13% of the proteome of this organism shares >30% sequence identity with BLAST hits against Firmicutes (img.jgi.doe.gov) so the BMC locus is likely part of a large HGT from a Firmicute. The closest relative to the BMC in this organism according to locus scoring is a EUT2I from *Clostridium scatologenes*, indicating a potential origin. In the case of RMM2 there are two phylum outliers, one in Epsilonproteobacteria and one in *Candidatus* Marinimicrobia (Supplementary Data 3). Both have phage integrase proteins (pfam13356 and pfam00589 domains) right next to the BMC locus, indicating a likely transfer via a phage vector.

Satellite loci. We define satellite loci as loci distal from the main locus and containing either shell proteins only, or a combination of shell proteins and enzymes. Some satellite loci seem to be obligately distal from the main BMC locus, such as the CcmK3/K4 paralogs of the beta-carboxysome;^{35,36} separate regulation of its expression may serve as a means to tune shell permeability under changing environmental conditions³⁵. Alteration of shell permeability may be a general function for shell proteins encoded in satellite loci of many types of BMCs (Fig. 2). Other satellite loci appear to have arisen as fissions from the main locus, such as the satellite locus for type I HO BMCs that contain two BMC-P with an aldolase, which are found in the main locus of type II HO BMCs (Supplementary Data 2). For HO, SPU3, and SPU6 BMCs some satellite loci resemble “EUT modules”, consisting of the ethanolamine degradation signature enzymes EutA/B/C and a EutL type BMC-T^{SP} shell protein (Supplementary Fig. 4). Because those BMC types are not expected to primarily process ethanolamine, this could represent a functional extension of the main BMC to use ethanolamine as an alternate substrate, with the BMC-T^{SP} acting as a shell protein that facilitates entry and the EutA/B/C enzymes to process it. In the case of SPU6 there is even an indication of integration of the satellite locus into the main locus, replacing the SPU type signature enzymes and the resulting locus is similar to that of EUT3 (Supplementary Fig. 4). EUT3 is a new type of ethanolamine utilization BMC with an AldDh that is phylogenetically distinct from the ones from both EUT1 and EUT2 loci (Fig. 5) and it contains a BMC-T^{DP} shell protein, unlike any other known EUT type BMC (Fig. 3d). Phylogenetic trees validate the link between the two locus types; they are on the same major branch of the AldDh tree (Fig. 5) and shell proteins like the BMC-T^{DP} are also adjacent, with the members of the fused locus found at the base of the SPU6 part of the branch (Fig. 3d), another example of how this expanded survey plausibly recounts BMC evolutionary history.

Shell protein trees and BMC identification. From all the BMC loci we have collected more than 40,000 shell protein sequences. Almost 19,000 of them are unique, highlighting their diversity despite a common function to form BMC shells. Among them are about 4900 BMC-P, 8000 BMC-H, 1550 BMC-HP, 1700 BMC-T^S, 1600 BMC-T^{SP}, and 1000 BMC-T^{DP}. A further reduced set of

those was used to build phylogenetic trees that illustrate their diversity (Fig. 3).

The BMC-P proteins resolve into five major clades: depicted in green, blue, purple, gray, and orange (Fig. 3a). Representatives of gray and orange clades always co-occur as two of three members of a BMC-P triplet. The third member of a given triplet varies across locus types, and is drawn from the green, blue or purple clades (Supplementary Fig. 3). The BMC-P proteins from these three clades frequently occur as the sole BMC-P gene in loci that lack BMC-P paralogs. The only exception is the alpha-carboxysomal CsoS4A and CsoS4B, BMC-P proteins that form a long stem of the purple clade; they always co-occur as pair. The multiplicity of BMC-P paralogs is unexpected because only 12 pentamers are needed to cap polyhedral structures and hints at an additional functional role for BMC-P proteins.

The defining, conserved primary structure of the BMC-H proteins, the pfam00936 domain is about 80 amino acids in length. However, extensions of up to 200 residues are observed, most frequently at the C-terminus and those BMC-H cluster on the phylogenetic tree (Fig. 3b). Each BMC locus type contains at least one BMC-H found close to the base of the tree (Fig. 3b). We refer to these as “basal” BMC-H proteins, because they share sequence motifs including the highly conserved inter-hexamer interface residue motifs KAAN and (P/A)RPH of the facets¹². Those likely constitute the bulk of the BMC shell facets while other BMC-H have more specialized functions.

Before the availability of genomic data, BMC-T^{DP} proteins were thought to be a rarity because their occurrence was limited to carboxysomal members CcmP and CsoS1D, as the two classic metabolosome model systems, EUT1 and PDU1, lack these proteins. In our analysis we find BMC-T^{DP} proteins in a large variety of BMC types (Fig. 2a). The BMC-T^{DP} tree can be divided into four major clades (Fig. 3d). The RMM1 and RMM2 BMC types are found in the clades shown in yellow. An unusual outlier is found at the base of the tree with homologs found in PDU1D loci from Proteobacteria and a single Acidobacterium; no other PDU type BMCs contain BMC-T^{DP} shell proteins. The major clade shown in blue contains a variety of uncharacterized BMCs (MIC12/BUF23) as well as EUT3 and SPU6 members that are in adjacent clades. The largest number of sequences are found in two clades that can be characterized by the presence of carboxysomal members. The clade depicted in purple contains the alpha-carboxysomal CsoS1D and a number of CsoS1D-like proteins from mainly HO and SPU3 and SPU5 type loci. The green-colored major clade contains the beta-carboxysomal CcmP as well as CcmP-like proteins from HO, MIC3, SPU2, and SPU3 type loci. The proximity of the BMC-T^{DP} of SPU to the carboxysome representatives suggests parallels with regard to the molecules that enter or leave the shell through these proteins. Sugar phosphates seem a likely candidate that both BMC types have in common, however BMC-T^{DP} from the other two major clades do not all share that type of substrate. It is possible that those BMCs have adopted the gated shell proteins for other purposes, such as large cofactors. Universally conserved across all types are the residues for the gating mechanism, indicating that this is a crucial function of the BMC-T^{DP}.

Shell proteins are diagnostic of the potential to form BMCs, in contrast to the enzymes found in BMCs, which have homologs that are not BMC-related. Scoring a proteome with the collection of BMC-type specific shell protein HMMs provides a quick assessment of both the presence of BMCs and initial identification due to their specificity, predicted from the color-type combinations of shell proteins present (Fig. 3 and Supplementary Data 2). We made a preliminary survey for the prevalence of BMCs in metagenomes. The shell protein HMMs derived from our locus collection were used to score data from 26,948 metagenomes in

IMG/M¹⁸, finding hits in 15,604 of them, and the total shell protein gene count adds up to more than 1.7 M. Using the assumption that the co-occurrence of the type-specific BMC-H and BMC-P colors are indicative of BMC functional type (or types that use the same or closely related substrates) we can make initial predictions the kinds of encapsulated metabolism to be found in metagenomes/microbiomes (Supplementary Fig. 5). The distribution of BMC types across bacterial clades, such as the prominent occurrence of PVM and SPU, provides a diagnostic marker of the specialized metabolism required in diverse nutritional landscapes of environmental samples. For example, the Anaerolineae class of Chloroflexi that contains SPU6 BMCs are prominent in environments characterized by anaerobic degradation of organic matter³⁷ and the Planctomycetes harboring PVM BMCs, despite growing very slowly, are one of the dominant bacterial species in algal blooms³⁸.

Discussion

By compiling over 40,000 BMC shell protein genes and surveying their genomic context, we find that the number of identifiable metabolic organelles has, over the last seven years, increased 20-fold, and representatives are now found across 45 phyla. The phylogenetic classification of shell proteins combined with the BMC type assignment shows that describing BMC shell proteins based on a specific locus type (e.g. PDU) is of limited usefulness as the growing number of functionally distinct BMC loci encode cohorts of shell proteins from across the phylogenetic tree, with many appearing in more than one locus type. Our phylogenetic classification also enables us to predict the basal hexamer(s) for a given locus (Fig. 3b), likely to form the bulk of the BMC shell facets. As such, it presumably conducts key metabolites; in conjunction with structural modeling focused on pore size and charge these data can be useful for predicting the first substrate of the encapsulated chemistry³⁹. In addition, we predict that basal BMC-H proteins are broadly interchangeable among functionally distinct BMCs, consistent with results of previous efforts to construct chimeric shells that, in retrospect, involved basal BMC-H proteins^{40,41}. The color combinations of shell proteins found in loci provide a guide to their compatibility in assembly and can be used for designing shells for metabolic engineering^{42–44}.

Our analysis also sheds light on the evolutionary history of BMCs. In addition to the widespread HGT of loci evident from their distribution across phyla (Fig. 4), individual shell proteins provide links, such as the BMC-T^{4p} of EUT3 and SPU6. Likewise, with our observations we can propose steps in the evolution of the encapsulated catalysis. For example the canonical PDU1AB and the recently described GRM4^{39,45}, all shell proteins have the same HMM derived color (Supplementary Data 3), indicating they are very closely related, and likely one of these BMC types originated from the other. Comparison of the locus diagrams reveals the same gene order, with the PDU signature enzymes interchanged with the GRM4 signature enzyme. Such a transformation, without altering the shell protein composition, is facilitated by sharing the initial substrate, 1,2-propanediol. This highlights the role of the shell proteins in shaping the encapsulated catalytic potential. Moreover, it suggests that finding the closest homologs from known BMCs, and the combination of shell protein color types, may be useful in predicting function in the absence of any information about the encapsulated enzymes, as in our preliminary survey of metagenomic data (Supplementary Fig. 5).

With the continued emphasis on sequencing using cultivation-independent methods which capture candidate phyla that are suggested to be the majority of bacteria⁴⁶, we expect that the role of BMCs in predicting the metabolic potential of environmental

samples will become even more prominent. For example, in this study, relative to the data available in 2014, we find the large expansion of the occurrence of SPU loci that is now represented by more than 280 loci across 26 bacterial phyla (eight members in 2014; Fig. 5). Their saprophytic function, and the availability of detritus in environmental samples, likely accounts for their emergence. We expect that with a continued focus on genome sequencing of the members of diverse ecosystems that the occurrence of FRAG loci will eventually be shown to correlate with an environmental factor that reveals the function of FRAG BMCs. These and other newly discovered functional types (Fig. 2) expand the types of predicted chemistry performed by BMCs and indicate that there is additional encapsulated dark biochemistry to be found.

Many microbial communities across Earth's biomes thrive in highly competitive environments, where their success to utilize resources and adapt under stable or changing environmental conditions will determine the fate of their genetic persistence. Our study finds BMCs in about 20% of all sequenced bacterial genomes. Their broad phylogenetic distribution across 45 bacterial phyla and wide environmental distribution spanning diverse ecosystems, underscores the important role they play in allowing bacteria to thrive in otherwise inaccessible environments. In addition to their prominence in candidate phyla and environmental samples, the importance of BMCs in organisms of the human microbiome²⁹ and their link to dysbiosis is also becoming apparent. For example, of the eleven bacterial species most prominently associated with complicated urinary tract infections⁴⁷, eight contain BMC loci, including PDU, EUT, GRM1, GRM2, GRM3, and FRAG; several of these commonly co-occur within a single genome. In the urinary tract, in addition to the breakdown of ethanolamine and 1,2-propanediol, the catabolism of choline via GRM1 and GRM2³⁹, likely serves as a carbon, nitrogen, and energy source that confers a competitive advantage to the uropathogens. More broadly, a species' potential to form multiple distinct BMCs functionally parallels the prevalence of BMCs in coexisting community members, foregrounds their previously undescribed role in providing metabolic flexibility as a driver for niche expansion. In addition to providing the foundation for understanding the native roles of BMCs in natural ecosystem function and dysbiosis, our catalog also provides insight into the diversity of metabolic compartmentalization and the evolutionary steps leading to this innovation that can inform engineering strategies.

Methods

Locus identification, clustering, and HMM generation. The complete Uniprot database was searched with shell protein pfams (PF03319 and PF00936), InterPro domain identifiers (IPR000249, IPR004992, IPR009193, IPR009307, IPR013501, IPR014076, IPR014077, IPR020808, IPR030983, IPR030984), PROSITE identifier PS01139, and SMART ID SM00877 [http://smart.embl-heidelberg.de/smart/do_annotation.pl?DOMAIN=SM00877], accessed as recently as March 2020. Adjacent proteins based on the numerical suffix were then obtained by downloading gene entries within 12 loci of the extracted gene identifier. Pfam tags were extracted from the uniprot annotations for initial HMM generation from proteins of the same pfam. HMMs were generated by aligning the sequences with clustalw 2.1⁴⁸, trimming with trimAl 1.2.rev59 with parameters -gt 0.6 -cons 30 -w 3⁴⁹ and HMMs built with hmmbuild from the HMMER package version 3.1b2⁵⁰. HMMs were calculated analogously for distinct shell proteins (see shell protein section below) and we then scored every locus protein against this combined HMM library, allowing us to represent each consolidated BMC locus as a string of identifier elements derived from the best-scoring HMM for each protein. To cluster the BMC loci and identify BMC types, we calculated pairwise correlation scores across all loci. Pairwise locus-locus scores were calculated amongst all loci using a python script based on the sum of two values: the total number of HMMs or pfams found in both loci, with each match multiplied by a weighting factor of 1.5 for each nearby (within three genes) shell protein gene; and the length of the longest sequence of consecutive HMM or pfam matches, multiplied by a weighting factor of 10. Data were then imported in Cytoscape 3.7.2⁵¹ to visualize the locus clustering. See Supplementary Fig. 6 for an all-vs-all BMC types visualization. Clusters

were manually matched to known BMC types and unknown types were assigned new identifiers. In a second round, the proteins that were not assigned a pfam in Uniprot and did not match another protein in the HMM library were collected and the whole set was then clustered with MMseqs2 (`mmseqs easy-cluster`⁵²). The clusters were then manually inspected and promising candidates were selected based on occurrence within a specific BMC type, distance to shell proteins, and direction of translation. The clusters were then analyzed and HMMs were generated from those proteins to identify them in the locus analysis. The overall prevalence of BMCs in genomes was assessed by scoring 64,495 isolate genomes, 1667 single-cell genomes (IMG/M¹⁸, April 2020) and 9,331 metagenome-assembled genomes (GEM¹⁹, February 2021) with a shell protein HMM library using `hmmsearch`⁵⁰ with a cutoff of 1E-20. The HMM collection can be found as a compressed file in Supplementary Data 5, and the HMM names are described in Supplementary Data 6.

Locus visualization. The HMM library was used to score each BMC type separately to generate type-specific HMMs using only the sequence from one type. This type-specific HMM was then used to score loci with `hmmsearch` and locus data were visualized using a python script (see Supplementary Data 3 for diagrams). The directionality of the gene coding for the proteins is shown as an arrow. While these data are not present in Uniprot files it can be determined by extracting the EMBL database identifier and parsing a corresponding DNA file downloaded through the European Nucleotide Archive (ENA). To determine the presence of EPs we collected sequences of known proteins with EPs and used the EP portion to generate HMMs. A separate HMM was generated for each BMC locus type and protein family. The combined HMM library was then used to identify potential EPs. Despite the low sequence conservation and short length of EPs this method is quite sensitive, however manual inspection of the results is still necessary.

Shell protein phylogenies. An initial set of 6408 BMC-P protein sequences was made non-redundant to 70% identity using the `usearch -cluster_fast` algorithm (v.11.0⁵³), resulting in a set of 1183 unique sequences. An initial alignment with `muscle`⁵⁴ was manually edited upon visual inspection with `Jalview`⁵⁵ to prune fragments and problematic sequences likely arising from genome sequencing or gene modeling errors. Sequences were then realigned using `MAFFT-linsi`⁵⁶ and uninformative columns were removed with `BMGE (-h 0.8 -g 0.05)`⁵⁷. `ModelTest-NG`⁵⁸ was used to determine the best-scoring substitution model, LG4M, which was then used to construct a maximum likelihood tree using `RAXML-NG (v0.6.0)`⁵⁹. A similar approach was used to collapse the redundancy of BMC-H proteins (95% identity), BMC-HP: (95%), BMC-T^s (90%), BMC-T^{dp} (90%), BMC-T^{sp} (90%); higher thresholds for these were necessary because of the higher overall homology. Trees were examined using `Archaeopteryx (www.phylosoft.org/archaeopteryx)` and significant clades were manually identified based on a general criterion of having a long internal stem. They were then assigned unique color names from the XKCD color survey (<https://xkcd.com/color/rgb/>) and colored with their corresponding RGB hexcodes, selecting similar colors for nearby clades. The sequences from each color-based clade were then subdivided by the initial BMC locus type assignments and used to generate subtype-specific HMMs for scoring the entire BMC dataset. The phylogenies were later re-examined with regards to final locus type assignments and annotated for functional correspondences to the clades or subclans (Fig. 3). Vector-quality images of the six phylogenies are provided as supplementary material with sequence identifiers and color assignments provided, for legibility upon manual zoom in a PDF reader (Supplementary Data 1) and additionally as an XML format file (Supplementary Data 4). `Pymol 2.3.0a0 (https://pymol.org)` was used to generate the structural figures in Supplementary Fig. 1.

Tree of bacterial phyla. A set of 56 universal single-copy marker proteins^{60,61} was used to build a phylogenetic tree for the domain Bacteria based on a representative dataset that included one genome for each bacterial order present in the IMG/M database (ref. 18, accessed March 2020). Marker proteins were identified with `hmmsearch (version 3.1b2)` using a specific HMM for each marker. Genomes lacking a substantial proportion of marker proteins (more than 26) or which had additional copies of more than five single-copy markers were removed from the dataset. For each marker, proteins were extracted, alignments built with `MAFFT-linsi (v7.294b)`⁶², and subsequently trimmed with `BMGE (v1.1257)` using `BLOSUM30`. Single protein alignments were then concatenated resulting in an alignment of 10,755 sites. Maximum likelihood phylogenies were inferred with `FastTree2`⁶³ using the options: `-spr 4 -mlacc 2 -slownni -wag`. The phylogenetic tree was pruned to keep only 1 representative genome for each phylum and then visualized using the `ete3` package⁶⁴.

Reporting summary. Further information on research design is available in the Nature Research Reporting Summary linked to this article.

Data availability

All sequence and metadata were obtained from public databases (Uniprot, ENA, IMG/M). All discussed BMC types are found as locus diagrams in Supplementary Data 3 and

the individual unique GeneIDs can be derived from the diagrams. An annotated version of the phylogenomic trees in Fig. 3 can be found as a pdf in Supplementary Data 1 as well as in XML format in Supplementary Data 4. The HMM collection can be found as a compressed file in Supplementary Data 5, the HMM names used in those are described in Supplementary Data 6. Any additional data is available from the corresponding author upon reasonable request.

Code availability

The Python code for scoring and visualizing BMC loci and comparing locus information is available at Github (https://github.com/markussutter/bmc_loci) and Zenodo (<https://doi.org/10.5281/zenodo.4794016>).

Received: 26 April 2021; Accepted: 28 May 2021;

Published online: 21 June 2021

References

- Drews, G. & Niklowitz, W. Beiträge zur Cytologie der Blaualgen. II. Zentroplasma und granuläre Einschlüsse von *Phormidium uncinatum*. *Arch. f. ür. Mikrobiologie* **24**, 147–162 (1956).
- Shively, J. M., Ball, F., Brown, D. H. & Saunders, R. E. Functional organelles in prokaryotes: polyhedral inclusions (carboxysomes) of *Thiobacillus neapolitanus*. *Science* **182**, 584–586 (1973).
- Shively, J. M. et al. Sequence homologs of the carboxysomal polypeptide CsoS1 of the thiobacilli are present in cyanobacteria and enteric bacteria that form carboxysomes - polyhedral bodies. *Can. J. Bot.* **76**, 906–916 (1998).
- Axen, S. D., Erbilgin, O. & Kerfeld, C. A. A taxonomy of bacterial microcompartment loci constructed by a novel scoring method. *PLoS Comput. Biol.* **10**, e1003898 (2014).
- Jorda, J., Lopez, D., Wheatley, N. M. & Yeates, T. O. Using comparative genomics to uncover new kinds of protein-based metabolic organelles in bacteria. *Protein Sci.* **22**, 179–195 (2013).
- Brinsmade, S. R., Paldon, T. & Escalante-Semerena, J. C. Minimal functions and physiological conditions required for growth of *Salmonella enterica* on ethanolamine in the absence of the metabolosome. *J. Bacteriol.* **187**, 8039–8046 (2005).
- Kerfeld, C. A. & Erbilgin, O. Bacterial microcompartments and the modular construction of microbial metabolism. *Trends Microbiol.* **23**, 22–34 (2015).
- Erbilgin, O., Sutter, M. & Kerfeld, C. A. The structural basis of coenzyme A recycling in a bacterial organelle. *PLoS Biol.* **14**, e1002399 (2016).
- Kerfeld, C. A., Aussignargues, C., Zarzycki, J., Cai, F. & Sutter, M. Bacterial microcompartments. *Nat. Rev. Microbiol.* **16**, 277–290 (2018).
- Kinney, J. N., Salmeen, A., Cai, F. & Kerfeld, C. A. Elucidating essential role of conserved carboxysomal protein CcmN reveals common feature of bacterial microcompartment assembly. *J. Biol. Chem.* **287**, 17729–17736 (2012).
- Aussignargues, C., Paasch, B. C., Gonzalez-Esquer, R., Erbilgin, O. & Kerfeld, C. A. Bacterial microcompartment assembly: the key role of encapsulation peptides. *Commun. Integr. Biol.* **8**, e1039755 (2015).
- Sutter, M., Greber, B., Aussignargues, C. & Kerfeld, C. A. Assembly principles and structure of a 6.5-MDa bacterial microcompartment shell. *Science* **356**, 1293–1297 (2017).
- Kalnins, G. et al. Encapsulation mechanisms and structural studies of GRM2 bacterial microcompartment particles. *Nat. Commun.* **11**, 388 (2020).
- Greber, B. J., Sutter, M. & Kerfeld, C. A. The plasticity of molecular interactions governs bacterial microcompartment shell assembly. *Structure* **27**, 749–763.e744 (2019).
- Sutter, M. et al. Structure of a synthetic beta-carboxysome shell. *Plant Physiol.* **181**, 1050–1058 (2019).
- Tanaka, S. et al. Atomic-level models of the bacterial carboxysome shell. *Science* **319**, 1083–1086 (2008).
- Kerfeld, C. A. et al. Protein structures forming the shell of primitive bacterial organelles. *Science* **309**, 936–938 (2005).
- Chen, I. A. et al. The IMG/M data management and analysis system v.6.0: new tools and advanced capabilities. *Nucleic Acids Res.* **49**, D751–D763 (2021).
- Nayfach, S. et al. A genomic catalog of Earth's microbiomes. *Nat. Biotechnol.* <https://doi.org/10.1038/s41587-020-0718-6> (2020).
- Sjolander, K. Phylogenomic inference of protein molecular function: advances and challenges. *Bioinformatics* **20**, 170–179 (2004).
- Krogh, A., Brown, M., Mian, I. S., Sjolander, K. & Haussler, D. Hidden Markov models in computational biology. Applications to protein modeling. *J. Mol. Biol.* **235**, 1501–1531 (1994).
- Waterhouse, A. et al. SWISS-MODEL: homology modelling of protein structures and complexes. *Nucleic Acids Res.* **46**, W296–W303 (2018).

23. Zhang, R. G. et al. The 2.2 Å resolution structure of RpiB/AlsB from *Escherichia coli* illustrates a new approach to the ribose-5-phosphate isomerase reaction. *J. Mol. Biol.* **332**, 1083–1094 (2003).
24. Aussenargues, C. et al. Structure and function of a bacterial microcompartment shell protein engineered to bind a [4Fe-4S] cluster. *J. Am. Chem. Soc.* **138**, 5262–5270 (2016).
25. Mohr, K. I. Diversity of myxobacteria—we only see the tip of the iceberg. *Microorganisms* <https://doi.org/10.3390/microorganisms6030084> (2018).
26. Erbilgin, O., McDonald, K. L. & Kerfeld, C. A. Characterization of a planctomycetal organelle: a novel bacterial microcompartment for the aerobic degradation of plant saccharides. *Appl. Environ. Microbiol.* **80**, 2193–2205, (2014).
27. Sizikov, S. et al. Characterization of sponge-associated Verrucomicrobia: microcompartment-based sugar utilization and enhanced toxin-antitoxin modules as features of host-associated Opitutales. *Environ. Microbiol.* **22**, 4669–4688 (2020).
28. Sichert, A. et al. Verrucomicrobia use hundreds of enzymes to digest the algal polysaccharide fucoidan. *Nat. Microbiol.* <https://doi.org/10.1038/s41564-020-0720-2> (2020).
29. Asija, K., Sutter, M. & Kerfeld, C. A. A survey of bacterial microcompartment distribution in the human microbiome. *Front. Microbiol.* <https://doi.org/10.3389/fmicb.2021.669024> (2021).
30. Losey, N. A. et al. *Thermoanaerobaculum aquaticum* gen. nov., sp. nov., the first cultivated member of Acidobacteria subdivision 23, isolated from a hot spring. *Int. J. Syst. Evol. Microb.* **63**, 4149–4157 (2013).
31. Fierer, N. Embracing the unknown: disentangling the complexities of the soil microbiome. *Nat. Rev. Microbiol.* **15**, 579–590 (2017).
32. Sturms, R., Streaulsin, N. A., Cheng, S. & Bobik, T. A. In *Salmonella enterica*, Ethanolamine utilization is repressed by 1,2-propanediol to prevent detrimental mixing of components of two different bacterial microcompartments. *J. Bacteriol.* **197**, 2412–2421 (2015).
33. Fendrich, C., Hippe, H. & Gottschalk, G. *Clostridium-Halophilium* Sp-Nov and *Clostridium-Litorale* Sp-Nov, an obligate halophilic and a marine species degrading betaine in the stickland reaction. *Arch. Microbiol.* **154**, 127–132 (1990).
34. Flint, H. J., Scott, K. P., Duncan, S. H., Louis, P. & Forano, E. Microbial degradation of complex carbohydrates in the gut. *Gut Microbes* **3**, 289–306 (2012).
35. Sommer, M. et al. Heterohexamers formed by CcmK3 and CcmK4 increase the complexity of beta carboxysome shells. *Plant Physiol.* **179**, 156 (2019).
36. Sommer, M., Cai, F., Melnicki, M. & Kerfeld, C. A. beta-Carboxysome bioinformatics: identification and evolution of new bacterial microcompartment protein gene classes and core locus constraints. *J. Exp. Bot.* **68**, 3841–3855 (2017).
37. Xia, Y., Wang, Y. B., Wang, Y., Chin, F. Y. L. & Zhang, T. Cellular adhesiveness and cellulolytic capacity in Anaerolineae revealed by omics-based genome interpretation. *Biotechnol. Biofuels* <https://doi.org/10.1186/s13068-016-0524-z> (2016).
38. Cai, F., Bernstein, S. L., Wilson, S. C. & Kerfeld, C. A. Production and characterization of synthetic carboxysome shells with incorporated luminal proteins. *Plant Physiol.* **170**, 1868–1877 (2016).
39. Zarzycki, J., Erbilgin, O. & Kerfeld, C. A. Bioinformatic characterization of glycol radical enzyme-associated bacterial microcompartments. *Appl. Environ. Microbiol.* **81**, 8315–8329 (2015).
40. Cai, F., Sutter, M., Bernstein, S. L., Kinney, J. N. & Kerfeld, C. A. Engineering bacterial microcompartment shells: chimeric shell proteins and chimeric carboxysome shells. *ACS Synth. Biol.* **4**, 444–453 (2015).
41. Slinger Lee, M. F., Jakobson, C. M. & Tullman-Ereck, D. Evidence for improved encapsulated pathway behavior in a bacterial microcompartment through shell protein engineering. *ACS Synthetic Biol.* <https://doi.org/10.1021/acssynbio.7b00042> (2017).
42. Frank, S., Lawrence, A. D., Prentice, M. B. & Warren, M. J. Bacterial microcompartments moving into a synthetic biological world. *J. Biotechnol.* **163**, 273–279 (2013).
43. Kerfeld, C. A. & Sutter, M. Engineered bacterial microcompartments: apps for programming metabolism. *Curr. Opin. Biotech.* **65**, 225–232 (2020).
44. Kirst, H. & Kerfeld, C. A. Bacterial microcompartments: catalysis-enhancing metabolic modules for next generation metabolic and biomedical engineering. *BMC Biol.* <https://doi.org/10.1186/s12915-019-0691-z> (2019).
45. Ferlez, B., Sutter, M. & Kerfeld, C. A. Glycol radical enzyme-associated microcompartments: redox-replete bacterial organelles. *mBio* <https://doi.org/10.1128/mBio.02327-18> (2019).
46. Castelle, C. J. & Banfield, J. F. Major new microbial groups expand diversity and alter our understanding of the tree of life. *Cell* **172**, 1181–1197 (2018).
47. Flores-Mireles, A. L., Walker, J. N., Caparon, M. & Hultgren, S. J. Urinary tract infections: epidemiology, mechanisms of infection and treatment options. *Nat. Rev. Microbiol.* **13**, 269–284 (2015).
48. Larkin, M. A. et al. Clustal W and clustal X version 2.0. *Bioinformatics* **23**, 2947–2948 (2007).
49. Capella-Gutierrez, S., Silla-Martinez, J. M. & Gabaldon, T. trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972–1973 (2009).
50. Eddy, S. R. Accelerated profile HMM searches. *PLoS Comput. Biol.* <https://doi.org/10.1371/journal.pcbi.1002195> (2011).
51. Shannon, P. et al. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res.* **13**, 2498–2504 (2003).
52. Steinegger, M. & Soding, J. MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nat. Biotechnol.* **35**, 1026–1028 (2017).
53. Edgar, R. C. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics* **26**, 2460–2461 (2010).
54. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* **32**, 1792–1797 (2004).
55. Waterhouse, A. M., Procter, J. B., Martin, D. M. A., Clamp, M. & Barton, G. J. Jalview Version 2—a multiple sequence alignment editor and analysis workbench. *Bioinformatics* **25**, 1189–1191 (2009).
56. Yamada, K. D., Tomii, K. & Katoh, K. Application of the MAFFT sequence alignment program to large data-reexamination of the usefulness of chained guide trees. *Bioinformatics* **32**, 3246–3251 (2016).
57. Criscuolo, A. & Gribaldo, S. BMGE (Block Mapping and Gathering with Entropy): a new software for selection of phylogenetic informative regions from multiple sequence alignments. *BMC Evol. Biol.* **10**, 210 (2010).
58. Darriba, D. et al. ModelTest-NG: a new and scalable tool for the selection of DNA and protein evolutionary models. *Mol. Biol. Evol.* **37**, 291–294 (2020).
59. Kozlov, A. M., Darriba, D., Flouri, T., Morel, B. & Stamatakis, A. RAXML-NG: a fast, scalable and user-friendly tool for maximum likelihood phylogenetic inference. *Bioinformatics* **35**, 4453–4455 (2019).
60. Yu, F. B. et al. Microfluidic-based mini-metagenomics enables discovery of novel microbial lineages from complex environmental samples. *Elife* <https://doi.org/10.7554/eLife.26580> (2017).
61. Eloe-Fadrosh, E. A. et al. Global metagenomic survey reveals a new bacterial candidate phylum in geothermal springs. *Nat. Commun.* **7**, 10476 (2016).
62. Katoh, K. & Standley, D. M. A simple method to control over-alignment in the MAFFT multiple sequence alignment program. *Bioinformatics* **32**, 1933–1942 (2016).
63. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2—approximately maximum-likelihood trees for large alignments. *PLoS ONE* **5**, e9490 (2010).
64. Huerta-Cepas, J., Serra, F. & Bork, P. ETE 3: reconstruction, analysis, and visualization of phylogenomic data. *Mol. Biol. Evol.* **33**, 1635–1638 (2016).
65. Klein, M. G. et al. Identification and structural analysis of a novel carboxysome shell protein with implications for metabolite transport. *J. Mol. Biol.* **392**, 319–333 (2009).
66. Cai, F. et al. The structure of CcmP, a tandem bacterial microcompartment domain protein from the beta-carboxysome, forms a subcompartment within a microcompartment. *J. Biol. Chem.* **288**, 16055–16063 (2013).
67. Larsson, A. M., Hasse, D., Valegard, K. & Andersson, I. Crystal structures of beta-carboxysome shell protein CcmP: ligand binding correlates with the closed or open central pore. *J. Exp. Bot.* **68**, 3857–3867 (2017).

Acknowledgements

This work was supported by the National Institutes of Health, National Institute of Allergy and Infectious Diseases (NIAID) grant 1R01AI114975-01 and the U.S. Department of Energy, Basic Energy Sciences, Contract DE-FG02-91ER20021. We thank Henning Kirst, Jan Zarzycki, and Stephanie A. Eichorst for helpful discussions, and Emiley Eloe-Fadrosh and Neha Varghese for custom HMM searches against the IMG/M database. Parts of this study were performed by the US Department of Energy Joint Genome Institute, a DOE Office of Science User Facility, under Contract No. DE-AC02-05CH11231 and made use of resources of the National Energy Research Scientific Computing Center, which is supported by the Office of Science of the US Department of Energy under Contract No. DE-AC02-05CH11231.

Author contributions

M.S. and C.A.K. conceived and designed the experiments. M.S., M.R.M., and F.S. performed the experiments. M.S., C.A.K., and T.W. analyzed the data and interpreted the results. M.S. and C.A.K. wrote the paper.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41467-021-24126-4>.

Correspondence and requests for materials should be addressed to C.A.K.

Peer review information *Nature Communications* thanks the anonymous reviewers for their contribution to the peer review of this work.

Reprints and permission information is available at <http://www.nature.com/reprints>

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021