

1 ***Title***

2 **A most wanted list of conserved protein families with no known domains**

3

4 ***Authors***

5 Stacia K. Wyman<sup>1,2</sup>, Aram Avila-Herrera<sup>1,3</sup>, Stephen Nayfach<sup>1,4,5</sup>, Katherine S. Pollard<sup>1,4,\*</sup>

6

7 ***Affiliations***

8 1: Gladstone Institutes, San Francisco, CA

9 2: University of California, Berkeley, CA

10 3: Lawrence Livermore National Laboratory, Livermore, CA

11 4: University of California, San Francisco, CA

12 5: DOE Joint Genome Institute, Walnut Creek, CA

13

14 \*: Corresponding Author ([kpollard@gladstone.ucsf.edu](mailto:kpollard@gladstone.ucsf.edu))

15

16 ***Abstract***

17 **The number and proportion of genes with no known function are growing rapidly. To**  
18 **quantify this phenomenon and provide criteria for prioritizing genes for functional**  
19 **characterization, we developed a bioinformatics pipeline that identifies robustly defined**  
20 **protein families with no annotated domains, ranks these with respect to phylogenetic**  
21 **breadth, and identifies them in metagenomics data. We applied this approach to 271 965**  
22 **protein families from the SFams database and discovered many with no functional**  
23 **annotation, including >118 000 families lacking any known protein domain. From these, we**  
24 **prioritized 6 668 conserved protein families with at least three sequences from organisms in**  
25 **at least two distinct classes. These Function Unknown Families (FUnkFams) are present in**  
26 **Tara Oceans Expedition and Human Microbiome Project metagenomes, with distributions**  
27 **associated with sampling environment. Our findings highlight the extent of functional**  
28 **novelty in sequence databases and establish an approach for creating a “most wanted” list of**  
29 **genes to characterize.**

30  
31 Genome sequencing and metagenomics are producing unprecedented amounts of data. But  
32 elucidation of gene function has not kept pace with the volume of identified genes. Homology-based  
33 annotation methods predict domains and functions for many new protein coding and RNA genes.  
34 However, many sequenced genes do not have significant homology to experimentally characterized  
35 domains or gene families. To quantify this problem, we developed a bioinformatics approach to  
36 identify *bona fide* protein families with no annotation and then characterized these with respect to  
37 their phylogenetic range and abundance in metagenomes. The result is FUnkFams, a prioritized  
38 catalog of genes for experimental discovery of function.

39

40 Our pipeline begins with a database of gene families, filters out truncated sequences without a start  
41 and stop codon, assigns annotations to all sequences in each family using one or more annotation  
42 databases, and records the taxonomy of the organism from which each sequence derived (Sup Fig  
43 S1). In a second step, metagenomic sequencing reads from a user-defined collection of samples are  
44 mapped to protein families, resulting in an estimate of protein family abundance in each sample.  
45 These data are then used to organize and rank gene families based on their level of annotation,  
46 number of sequences, phylogenetic diversity, and distribution across metagenomes.

47  
48 We applied this approach to discover the least annotated, most phylogenetically diverse full-length  
49 protein families in the SFams database (Sharpton et al., 2012) (Supplemental Text). We used  
50 SFams, because it was compiled in a comprehensive, automated fashion from thousands of diverse  
51 genome sequences, and we applied bioinformatics filters to remove small and truncated families.  
52 Specifically, we first identified 224 409 SFams with at least three unique, homologous, full-length  
53 protein sequences. We then annotated the sequences in these SFams using two curated and  
54 frequently updated sources of protein domains: the PFam database (Finn et al., 2014) and the NCBI  
55 Conserved Domain Database (CDD) (Marchler-Bauer et al., 2011). This analysis showed that the  
56 majority of protein families lack even a single domain annotated in PFam or CDD (N=118 607  
57 SFams, 52.9% of total). These protein families without domain annotation are comprised of  
58 sequences from many branches of the cellular tree of life (Fig 1A). For further analysis and top  
59 prioritization we selected a subset of 6 668 protein families with no annotated domains and  
60 sequences from two or more taxonomic classes (Sup Fig 2A). We call these Function Unknown  
61 Families (FUnkFams)(Sup Table S1). Most FUnkFams (84.3%) are not in UniProt xref (UniProt,  
62 2015), and those that are in UniProt are largely annotated as hypothetical or uncharacterized  
63 proteins (Sup Table S2).

64

65 FUnkFams are similar to other SFams in terms of properties other than the criteria we used to  
66 define them (i.e., functional annotation and phylogenetic breadth). Protein sequences in FUnkFams  
67 have a similar phylogenetic distribution to all SFams (Sup Fig S2B-C) with some enrichment in  
68 Cyanobacteria. They are also somewhat depleted in eukaryotes and archaea, probably due to  
69 bacterial SFams being more likely to meet our criteria of multiple homologous sequences from at  
70 least two classes. Like SFams, a typical FUnkFam is approximately 250 amino acids long (Fig 1B)  
71 and is comprised of three to five sequences (Fig 1C), though FUnkFams are slightly depleted for  
72 very long and very large families compared to better-annotated SFams. Nonetheless, six FUnkFams  
73 are comprised of more than 100 sequences, including a Proteobacterial family (SFams.ID=4560)  
74 with 203 sequences and a family (SFams.ID=5980) with 145 sequences spanning multiple domains  
75 of life. Thus, FUnkFams appear to be representative of full-length, phylogenetically diverse protein  
76 families.

77  
78 To investigate the ecological distributions of FUnkFams, we quantified their abundance in shotgun  
79 metagenomes from the Tara Oceans Expedition (TO; 243 samples from 210 ecosystems in 20  
80 biogeographic provinces at different depths over the course of three years) (Pesant et al., 2015) and  
81 Human Microbiome Project (HMP; 699 samples from oral, airways, skin, gut, vaginal sites on 300  
82 healthy individuals at up to three time points over two years) (Human Microbiome Project, 2012)  
83 (Supplemental Methods). The majority of FUnkFams (56.6%) are present in at least one of these  
84 942 metagenomes, with many detected in multiple metagenomes (32.5% in at least two HMP  
85 samples, 37.2% in at least two TO samples) but relatively few (13.3%) detected in both TO and  
86 HMP (Fig 2A). FUnkFam prevalence was generally higher in TO (mean=18.6% versus 8.1%), with  
87 TO samples averaging 700 detected FUnkFams and HMP averaging 304 (Sup Fig S3A-B). Higher  
88 sequencing depth in TO may contribute to this signal. Abundance of detected FUnkFams is on  
89 average higher in TO, though many FUnkFams are approximately equally abundant between TO

90 and HMP (Fig 2B; Sup Fig S3C) and 27 are highly abundant in both environments (Sup Fig S4).  
91 Reflecting the ecological specificity of many FUnkFams, beta-diversity is significantly higher  
92 between the two environments than between samples within either environment (Fig 2C).  
93  
94 We next used logistic regression to quantify how these differences in FUnkFam distributions across  
95 TO and HMP correlate with characteristics of the samples after adjusting for technical variables  
96 (Supplemental Methods). In TO, the presence of three FUnkFams was significantly associated with  
97 nitrate level after multiple testing correction (FDR<5%). One of these FUnkFams was also  
98 significantly associated with salinity and longitude, and another was significantly associated with  
99 longitude, latitude, temperature, and depth (Sup Table S3). Other FUnkFams showed weaker  
100 associations with environmental variables (Sup Fig S5, S6). The dominant variable associated with  
101 FUnkFam presence in HMP samples is body site (Sup Fig S5, S7; Sup Table S4), with only a few  
102 FUnkFams broadly detected across body sites. Other host phenotypes, such as BMI, smoking status,  
103 or diet, were not significantly associated with the presence of any FUnkFams.  
104  
105 These results identify thousands of uncharacterized protein families composed of homologous  
106 sequences from phylogenetically diverse organisms that are abundant in the human body or global  
107 oceans. These characteristics suggest that FUnkFams are *bona fide* protein families, and the  
108 associations of specific FUnkFams with marine environments or body sites provide hints about  
109 protein function and ecology. This study therefore lays the groundwork for significant future work  
110 to (i) predict (e.g., via genome proximity and further metagenome profiling (Lobb and Doxey, 2016)  
111 or literature based similarity (Price and Arkin, 2017)) and (ii) experimentally validate (e.g., via  
112 biochemical and structural characterization (Gerlt et al., 2011)) the functions of FUnkFams and the  
113 unannotated protein domains they contain. Our approach can be flexibly extended to use other

114 databases of gene families and sources of functional annotation, and it will be interesting to apply it  
115 to other protein catalogs as well as RNA genes.

116

117 Supplementary information is available at ISME Journal's website. FUnkFams data are freely  
118 available via figshare at:

119 [https://figshare.com/projects/Function\\_Unknown\\_Families\\_of\\_homologous\\_proteins\\_FUnkFams\\_/](https://figshare.com/projects/Function_Unknown_Families_of_homologous_proteins_FUnkFams_/25924)

120 25924.

121

## 122 **Acknowledgements**

123 This project was funded by NSF (#DMS-1563159), Gordon & Betty Moore Foundation (#3300), and  
124 the Gladstone Institutes. We thank Thomas Sharpton and Jonathan Eisen for very helpful  
125 conversations at the conception of the project.

## 126 References

- 127 BUCHFINK, B., XIE, C. & HUSON, D. H. 2015. Fast and sensitive protein alignment using  
128 DIAMOND. *Nat Methods*, 12, 59-60.
- 129 FINN, R. D., BATEMAN, A., CLEMENTS, J., COGGILL, P., EBERHARDT, R. Y., EDDY, S. R.,  
130 HEGER, A., HETHERINGTON, K., HOLM, L., MISTRY, J., SONNHAMMER, E. L., TATE, J.  
131 & PUNTA, M. 2014. Pfam: the protein families database. *Nucleic Acids Res*, 42, D222-  
132 30.
- 133 FOSTER, Z. S., SHARPTON, T. J. & GRUNWALD, N. J. 2017. Metacoder: An R package for  
134 visualization and manipulation of community taxonomic diversity data. *PLoS*  
135 *Comput Biol*, 13, e1005404.
- 136 GERLT, J. A., ALLEN, K. N., ALMO, S. C., ARMSTRONG, R. N., BABBITT, P. C., CRONAN, J. E.,  
137 DUNAWAY-MARIANO, D., IMKER, H. J., JACOBSON, M. P., MINOR, W., POULTER, C. D.,  
138 RAUSHEL, F. M., SALI, A., SHOICHET, B. K. & SWEEDLER, J. V. 2011. The Enzyme  
139 Function Initiative. *Biochemistry*, 50, 9950-62.
- 140 HUMAN MICROBIOME PROJECT, C. 2012. Structure, function and diversity of the healthy  
141 human microbiome. *Nature*, 486, 207-14.
- 142 LOBB, B. & DOXEY, A. C. 2016. Novel function discovery through sequence and structural  
143 data mining. *Curr Opin Struct Biol*, 38, 53-61.
- 144 MARCHLER-BAUER, A., LU, S., ANDERSON, J. B., CHITSAZ, F., DERBYSHIRE, M. K.,  
145 DEWEESE-SCOTT, C., FONG, J. H., GEER, L. Y., GEER, R. C., GONZALES, N. R., GWADZ,  
146 M., HURWITZ, D. I., JACKSON, J. D., KE, Z., LANCZYCKI, C. J., LU, F., MARCHLER, G. H.,  
147 MULLOKANDOV, M., OMELCHENKO, M. V., ROBERTSON, C. L., SONG, J. S., THANKI,  
148 N., YAMASHITA, R. A., ZHANG, D., ZHANG, N., ZHENG, C. & BRYANT, S. H. 2011. CDD:  
149 a Conserved Domain Database for the functional annotation of proteins. *Nucleic*  
150 *Acids Res*, 39, D225-9.
- 151 NAYFACH, S. & POLLARD, K. S. 2015. Average genome size estimation improves  
152 comparative metagenomics and sheds light on the functional ecology of the human  
153 microbiome. *Genome Biol*, 16, 51.
- 154 PESANT, S., NOT, F., PICHERAL, M., KANDELS-LEWIS, S., LE BESCOT, N., GORSKY, G.,  
155 IUDICONE, D., KARSENTI, E., SPEICH, S., TROUBLE, R., DIMIER, C., SEARSON, S. &  
156 TARA OCEANS CONSORTIUM, C. 2015. Open science resources for the discovery and  
157 analysis of Tara Oceans data. *Sci Data*, 2, 150023.
- 158 PRICE, M. N. & ARKIN, A. P. 2017. PaperBLAST: Text Mining Papers for Information about  
159 Homologs. *mSystems*, 2, e00039-17.
- 160 SHARPTON, T. J., JOSPIN, G., WU, D., LANGILLE, M. G., POLLARD, K. S. & EISEN, J. A. 2012.  
161 Sifting through genomes with iterative-sequence clustering produces a large,  
162 phylogenetically diverse protein-family resource. *BMC Bioinformatics*, 13, 264.
- 163 UNIPROT, C. 2015. UniProt: a hub for protein information. *Nucleic Acids Res*, 43, D204-12.
- 164

165

166 **Figure Captions:**

167 **Figure 1. A.** Phylogenetic heat tree of proteins in FUnkFams generated with Metacoder (Foster et  
168 al., 2017). Each FUnkFams protein sequence was annotated with the taxonomic label of the genome  
169 from which it was derived. The color of a branch represents the number of proteins from any  
170 FUnkFam on that branch of the taxonomy. The tree shows that FUnkFams are present across  
171 diverse lineages of cellular organisms including families from all three domains and over thirty  
172 phyla. Proteobacteria contribute many sequences to FUnkFams, in part because many genomes  
173 have been sequenced from that phylum. Supplemental Figure 2 shows the heat tree of all SFams,  
174 illustrating lineages where FUnkFams are enriched given how many genomes have been sequenced.

175 **B.** FUnkFams protein length (in amino acids, log<sub>2</sub> scale) and family size (number of protein  
176 sequences) are comparable to other SFams. Top and middle panels show histograms, and bottom  
177 panels are quantile-quantile plots showing that most quantiles of length and size are equal between  
178 FUnkFams and SFams, except at the top quantiles where SFams are slightly longer (i.e., more amino  
179 acids) and bigger (i.e., more sequences).

180 **Figure 2.** FUnkFams are present in marine and human metagenomes. **A.** Most FUnkFams are  
181 detected in either TO or HMP metagenomes, but relatively few are present in both environments. **B.**  
182 Heatmap showing the normalized abundance (RPKG) of FUnkFams (rows) in TO (left) or HMP  
183 (right) metagenomes. The 180 FUnkFams with at least 50 aligned reads across all samples are  
184 displayed (see Sup Fig S4 for the heatmap of all FUnkFams). **C.** Distributions of Bray-Curtis  
185 dissimilarity between pairs of samples from marine environments (TO; blue), between pairs of  
186 samples from human microbiomes (HMP; red), and between pairs of samples from different  
187 environments (white). Samples are more similar within than between the two environments.



## 188 **Supplemental Text**

### 189 **FUnkFams Construction**

190 Our pipeline (Sup Fig S1) begins with 345 641 SFam protein families that were previously derived  
191 from *de novo* iterative clustering of ~10.5 million protein sequences from ~3 000 diverse genomes  
192 (Sharpton et al., 2012). SFams with less than three unique protein sequences are dropped, and then  
193 SFams where >50% of the sequences lack a start or stop codon are removed. This rigorous filtering  
194 produced 224 409 full-length protein families, at the cost of eliminating some small SFams. We then  
195 searched for all the proteins in these full-length SFams in the PFam database (acquired via  
196 Swissprot) (Finn et al., 2014) and the NCBI Conserved Domain Database (CDD) (Marchler-Bauer et  
197 al., 2011) to annotate domains in every sequence. These database searches were to identify the  
198 exact protein from an SFam (100% identical blast hit over the full length of the SFam sequence), not  
199 to identify homologs of the SFams sequences. The rationale for this strategy is that the SFam  
200 sequences derive from genomes that have been processed into these databases, and hence any  
201 proteins from these genomes should have been annotated already based on homology and other  
202 criteria of the databases. We identified 118 607 families for which no sequences had any domain  
203 annotation in PFam or CDD. Next, we characterized each sequence in each protein family according  
204 to the NCBI taxonomic annotation of the genome from which it derived and then quantified how  
205 many different species, genera, families, orders, classes, phyla, kingdoms, and domains are  
206 represented in each gene family (Sup Fig S2). We found that from the 118 607 families with no  
207 annotated domains there are 6 668 families that contain sequences from at least two classes. We  
208 call these FUnkFams. We identified 1 045 FUnkFams with at least one sequence in UniProt's xref  
209 database (UniProt, 2015) (Sup Table S2). These are nearly all uncharacterized proteins, although  
210 we did identify eight FUnkFams with a sequence that has an xref-annotated function, despite having  
211 no domain annotation (Sup Table S5).

## 212 **Profiling in TO and HMP Metagenomes**

213 We used Diamond (Buchfink et al., 2015) to align each read in the TO and HMP metagenome  
214 samples to a database of SFams sequences. We counted aligned reads for each FUnkFam, requiring  
215 a best hit to a protein belonging to the FUnkFam with at least 99% DNA sequence identity over the  
216 whole length of the read. FUnkFams with at least one read count were called present in the  
217 metagenome. FUnkFam abundance in each sample was estimated using reads per kilobase of  
218 genome (RPKG), a statistic that normalizes for both protein family length (mean of all member  
219 sequences in SFams database) and average genome size (estimated from the metagenomics sample  
220 with MicrobeCensus) (Nayfach and Pollard, 2015).

221 The TO dataset had hits to 2 488 FUnkFams in 308 samples, and the HMP dataset had hits to 2 164  
222 FUnkFams in 696 samples (Fig 2A). Of these, 889 FUnkFams were detected in both environments,  
223 and these had a range of abundance levels across both TO and HMP (Sup Fig S3C). A particularly  
224 prevalent set of 137 FUnkFams was found in over 90% of TO samples, while just three were in over  
225 90% of HMP samples, likely reflecting greater annotation of functions found in the human body  
226 samples relative to marine samples but also potentially also due to ecological differences between  
227 human body sites.

## 228 **Association Testing in HMP**

229 We tested for association between a number of host phenotypes and FUnkFam presence in HMP  
230 metagenomes. To pre-filter FUnkFams without sufficient variation in presence across samples to  
231 detect associations, we only included FUnkFams with entropy in the top 25 percentile. To focus on  
232 the most phylogenetically diverse protein families, we additionally only included FUnkFams with  
233 sequences derived from genomes in at least two phyla (Sup Fig S2). This resulted in a set of 319  
234 FUnkFams for association testing. We investigated associations with 13 host phenotypes that  
235 reflect lifestyle and medication use, as defined in HMP documentation (Sup Table S6). Phenotype

236 data was obtained with permission through dbGaP (study ID = phs000228.v2.p1). Phenotypes were  
237 required to have at least two values with more than four observations. Seven subject variables  
238 passed this filtering step: bmi category, contraceptive use, breastfed status, diet, education level,  
239 birth country and student status. We fit a logistic regression model for each FUnkFam and used the  
240 resulting coefficients and their standard errors to perform t-tests to identify phenotypes associated  
241 with the presence of each FUnkFam across samples from each body site. The models account for  
242 geographic location (SITE variable in HMP) and were fit for each body site. P-values were corrected  
243 for multiple testing using the false discovery rate (FDR). We repeated this analysis within body  
244 subsites using the same filtering criteria and the resulting set of 335 FUnkFams (Sup Table S4).

245

#### 246 **Association Testing in Tara Oceans Data**

247 We tested for association between environmental variables and FUnkFam presence across TO  
248 samples. Environmental data was downloaded from the Tara Oceans data resource ([http://ocean-](http://ocean-microbiome.embl.de)  
249 [microbiome.embl.de](http://ocean-microbiome.embl.de)). Using the same criteria as with HMP, we analyzed only 100 FUnkFams with  
250 high entropy and sequences from at least two phyla. We fit logistic regression models for FUnkFam  
251 presence versus environmental variables, adjusting for latitude and month. Separate models were  
252 fit for samples collected with each filter size (size fraction). The resulting t-test p-values were  
253 adjusted for multiple testing using FDR (Sup Table S3).

254 **Supplemental Figures:**

255 **S1.** Bioinformatics pipeline for identifying FUnkFams from the SFams database

256 **S2.** A) Number of FUnkFams found across multiple domains, phyla, and classes in the tree of  
257 cellular organisms (e.g. 208 FUnkFams were found across more than one domain). B) Metacoder  
258 phylogenetic heat tree of SFams abundance across cellular organisms. Color indicates number of  
259 sequences on a branch. A random subset of 400 000 SFams was used to generate the tree. C)  
260 Metacoder phylogenetic heat tree of FUnkFams abundance across cellular organisms (as in Fig. 1A,  
261 for comparison here with SFams tree).

262 **S3.** A) Prevalence (vertical axis) of FUnkFams in TO (blue) and HMP (red) samples, ordered by  
263 decreasing prevalence in HMP (horizontal axis). B) Prevalence (vertical axis) of FUnkFams in TO  
264 (blue) and HMP (red) samples, ordered by decreasing prevalence in TO (horizontal axis). Many  
265 FUnkFams are more prevalent in TO than HMP, but the converse is not true. C) For 889 FUnkFams  
266 present in at least one TO and at least one HMP sample, the fractional abundance (vertical axis)  
267 represents the proportion of total RPKG for the FUnkFam that comes from TO (blue) versus HMP  
268 (red). FUnkFams are ordered by decreasing proportion of total RPKG deriving from TO samples  
269 (horizontal axis).

270 **S4.** Heatmap with all 3 763 FUnkFams (rows) detected in any metagenome (TO, HMP or both) at  
271 any abundance. Blue (left columns) are TO samples and red (right columns) are HMP samples.

272 **S5.** PCA plots of samples from HMP (A-B) and TO (C-E) based on counts of metagenomic sequencing  
273 reads mapped to all FUnkFams. HMP samples cluster by body site (A) but not other phenotypes  
274 such as BMI (B). TO samples cluster by marine layer (E) but not other environmental features (C-  
275 D).

276 **S6.** Heatmap for most abundant FUnkFams in TO samples, clustered both by column (samples) and  
277 row (FUnkFams) with environmental features annotated across rows.

278 **S7.** Heatmap for most abundant FUnkFams in HMP samples, clustered both by column (samples)  
279 and row (FUnkFams) with host phenotypes annotated across rows.

280

281 **Supplemental Tables:**

282 **Supplemental Table S1.** Characteristics of FUNkFams, including phylogenetic distribution and  
 283 prevalence in TO and HMP samples.

284 **Supplemental Table S2.** Annotations for 1 045 FUNkFams with a protein sequence in the UniProt  
 285 xref database.

286 **Supplemental Table S3.** Results of statistical tests for associations between environmental  
 287 variables and FUNkFams presence across TO samples.

288 **Supplemental Table S4.** Results of statistical tests for associations between host phenotype  
 289 variables and FUNkFams presence across HMP samples.

290 **Supplemental Table S5.** Annotations for eight FUNkFams with a protein sequence whose function  
 291 is annotated in the UniProt xref database (despite having no annotated domains).

292 **Supplemental Table S6.** HMP phenotypes tested for association with FUNkFams abundance.

293 BMI\_CAT has three levels: lean=BMI<25, overweight=BMI 25-30, obese=BMI30+.

HMP Variable	Description
DSUDIET	Meat/poultry frequency 294
DSUBFED	Breastfed as child
BRTHCTRY	Born in USA or Canada
EDLVL_BS	Holds BS degree or greater 295
OCPTN_ST	Is student
BMI_CAT	BMI categories
SMOKER	Uses tobacco and smokes cigarettes
DCMCODE_3	Subject was taking antacids during the visit
DCMCODE_8	Subject was taking antibiotics during the visit
DCMCODE_13	Subject was taking contraceptives during the visit
DCMCODE_16	Subject was taking GI meds during the visit
DCMCODE_18	Subject was taking hormones/steroids during the visit
DCMCODE_NA	Subject was not taking medication during the visit



Figure 2

