

# Pairing Metagenomics and Metaproteomics to Pinpoint Ecological Niches and Metabolic Essentiality of Microbial Communities

Tong Wang<sup>1,†</sup>, Leyuan Li<sup>2,3,†</sup>, Daniel Figeys<sup>3,\*</sup>, Yang-Yu Liu<sup>1,\*</sup>

<sup>1</sup>*Channing Division of Network Medicine, Department of Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, MA 02115, USA.*

<sup>2</sup>*State Key Laboratory of Proteomics, Beijing Proteome Research Center, National Center for Protein Sciences (Beijing), Beijing Institute of Lifeomics, Beijing 102206, China*

<sup>3</sup>*School of Pharmaceutical Sciences, Ottawa Institute of Systems Biology, Faculty of Medicine, University of Ottawa, Ottawa, ON K1H8M5, Canada.*

† These authors contributed equally to this work.

\* Correspondence: [dfigeys@uottawa.ca](mailto:dfigeys@uottawa.ca), [yyl@channing.harvard.edu](mailto:yyl@channing.harvard.edu)

## Abstract

Microbial genomes encode functional repertoire of microbes. However, microbes rely on various proteins to be expressed to carry out specific functions, and the expression of those proteins can be affected by the environment. It remains elusive how the selective expression of a protein depends on whether it is metabolically essential to the microbe's growth, or it can claim resources as an ecological niche. Here we show that by pairing metagenomics and metaproteomics data we can reveal whether a protein is relevant for occupying ecological niches or is essential for microbial metabolism. In particular, we developed a computational pipeline based on the quantification of the gene-level (or protein-level) functional redundancy of each protein, which measures the degree to which phylogenetically unrelated taxa can express (or have already expressed) the same protein, respectively. We validated this pipeline using both simulated data of a consumer-resource model and real data of human gut microbiome samples. Furthermore, for the real data, we showed that the metabolic and ecological roles of ABC-type transporters and ribosomal proteins predicted by our pipeline agree well with prior knowledge. Finally, we performed *in vitro* culture of a human gut microbiome sample and investigated how oversupplying various sugars involved in ecological niches influences the community structure and protein expression. The presented results help us identify metabolic and ecological roles of proteins, which will inform the design of nutrient interventions to modulate the human microbiome.

## Introduction

The advance in metagenomic sequencing technology has enabled us to measure the genomic contents and functional potentials of microbial communities at an unprecedented rate, helping us understand how the functionality of microbes influences host health<sup>1-3</sup> and how microbial metabolism in natural environments enables biogeochemical cycling<sup>4-6</sup>. Based on metabolic models inferred from genomes, various computational approaches have been proposed to better quantify inter-species interactions and ecological concepts in microbial communities<sup>7-12</sup>. For example, metabolic networks of microbes have been employed to quantify complementarity and competition indices as a proxy for potential interactions<sup>7</sup>. Also, a nonlinear dimensionality reduction technique has been used to map bacterial metabolic niche space<sup>9</sup>. In addition, functional redundancy and functional stability for microbial communities were analyzed in the past<sup>10-12</sup>.

A major limitation of those approaches is that they only rely on metagenomic data, which does not reflect true functional activities but only encodes functional capacity (or potential functions). In reality, at any given time and under any environmental condition, microbes only express a subset of their potential functions as proteins to carry out particular functions<sup>13</sup>. Recently, an ultra-deep metaproteomics approach has been developed to quantify expressed proteins in complex microbial communities, e.g., the human gut microbiome<sup>14</sup>. Pairing metagenomic and metaproteomic data offers the possibility to investigate how each protein is selectively expressed under different environmental conditions.

From the metabolic perspective, it is well known that some genes and their expressed proteins are indispensable for cell metabolism under any conditions, and microbes will not survive or reproduce if those genes are lost or those proteins are not expressed. Indeed, lacking proteins essential to microbial metabolism will cease microbial growth, regardless of ecological competition. For example, the growth of microbes relies on aminoacyl-tRNA<sup>15,16</sup>. Consequently, microbes have to express proteins involved in the aminoacyl-tRNA synthesis due to their metabolic essentiality to microbial growth<sup>15,16</sup>.

From the ecological perspective, some proteins are expressed under ecological selection, and the presence of such proteins directly indicates which resources a microbe can utilize so as to thrive, i.e., the ecological niche of this microbe in the microbial community. Different proteins might enable a microbe to utilize different resources or adapt to varying environments. If the function of a protein can simply be performed by another protein, it may

be not necessary to express both proteins at the same time. This is evident in the case of *E. coli*, which prefers glucose over lactose due to the repressed expression of lactose-utilizing enzymes, even though it can use both sugars<sup>17,18</sup>. Such specialization of consuming one resource caused by the selective protein expression may reduce the niche overlap with other species. Another example is Acetyl-coenzyme synthetase (Acs) --- a protein that catalyzes the conversion of acetate into Acetyl-CoA, an essential intermediate in the metabolism<sup>19,20</sup>. The overexpression of Acs in *E. coli* can significantly switch glucose consumption to acetate consumption<sup>21-24</sup>. The glucose specialist (CV103) and acetate specialist (CV101) are two *E. coli* mutants with different metabolic strategies; CV103 does not express Acs while CV101 overexpresses it<sup>21-24</sup>. It has been shown that CV101 can consume acetate produced by CV103, and thus they achieve a coexistence due to the niche partitioning<sup>21,22</sup>.

How to understand the selective expression of microbial proteins is an outstanding question in microbiology. Does the behavior of selective expression of microbial proteins differ between metabolic function (e.g., essential for microbial growth metabolism) and ecological function (e.g., claiming resources as a niche)? To address this question, in this work we developed a computational method to perform paired metagenomic and metaproteomic<sup>25-28,14</sup> data analysis and revealed whether a protein is essential for microbial metabolism or relevant for occupying ecological niches. In particular, we used the metagenomic data to construct the Gene Content Network (GCN) --- a bipartite graph that connects microbial taxa to their genes (Fig. 1a), and used the metaproteomic data to construct the Protein Content Network (PCN) -- a bipartite graph that connects microbial taxa to their truly expressed protein functions (Fig. 1b). For each protein, we quantified its gene-level (or protein-level) functional redundancy (FR), which is defined as the degree to which unrelated taxa can express (or have already expressed) this protein, respectively. Using synthetic data generated by a consumer-resource model of microbial communities, we found that either the comparison of network degree of a protein (i.e., the number of taxa that own/express the protein) between the GCN and PCN or the comparison between the gene-level and protein-level FR of a protein can reveal its role in metabolic essentiality and ecological niches. Then we applied the same computational pipeline to analyze the real data of human gut microbiome samples to predict metabolic and ecological functions for proteins. We found that the metabolic and ecological roles of ABC-type transporters and ribosomal proteins predicted by our method agree well with prior knowledge. Finally, we performed *in vitro* culture experiments using human gut microbiome samples with and without sugars added to investigate how oversupplying various sugars involved in ecological niches influences the community structure and protein expression.

## Results

### Quantifying gene- and protein-level functional redundancy of each protein

Consider a microbiome sample with taxonomic profile  $\mathbf{p} = (p_1, \dots, p_N)$ , where  $p_i$  is the relative abundance of taxon- $i$  and  $\sum_{i=1}^N p_i = 1$ . For a given protein, we can define its gene-level FR ( $\text{FR}_g$ ) within this sample as

$$\text{FR}_g = \sum_{i=1}^N \sum_{j \neq i}^N (1 - d_{ij}^{\text{GCN}}) p_i p_j, \quad (1)$$

where  $d_{ij}^{\text{GCN}}$  is the distance between taxon- $i$  and taxon- $j$  based on their genomic capacity to express this protein. For simplicity, we assume  $d_{ij}^{\text{GCN}}$  is binary, i.e.,  $d_{ij}^{\text{GCN}} = 0$  if and only if both taxa share the potential to express the protein, and  $d_{ij}^{\text{GCN}} = 1$  otherwise. For the same protein, we can also define its protein-level FR ( $\text{FR}_p$ ) within this sample as

$$\text{FR}_p = \sum_{i=1}^N \sum_{j \neq i}^N (1 - d_{ij}^{\text{PCN}}) p_i p_j, \quad (2)$$

where  $d_{ij}^{\text{PCN}}$  is the distance between taxon- $i$  and taxon- $j$  based on their expression of the protein. Again, we assume  $d_{ij}^{\text{PCN}}$  is binary, i.e.,  $d_{ij}^{\text{PCN}} = 0$  if and only if both taxa have expressed the protein, and  $d_{ij}^{\text{PCN}} = 1$  otherwise. Note that here we define  $\text{FR}_g$  and  $\text{FR}_p$  for each protein. This is different from our previous studies<sup>12,14</sup>, where FR was calculated by including all genes or proteins in the entire microbial community.

To demonstrate the definitions of  $\text{FR}_g$  and  $\text{FR}_p$ , let's consider a simple community consisting of two coexisting *E. coli* strains CV101 and CV103 with relative abundance  $p_1$  and  $p_2$ , respectively<sup>21,22</sup>. For the protein Acs that is required for the acetate consumption, since both CV101 and CV103 own this functional capacity, we have  $d_{12}^{\text{GCN}} = d_{21}^{\text{GCN}} = 0$ , and  $\text{FR}_g = 2p_1p_2$ . However, because CV103 does not express Acs and CV101 overexpresses it<sup>21-24</sup>, we have  $d_{12}^{\text{PCN}} = d_{21}^{\text{PCN}} = 1$ , and  $\text{FR}_p = 0$ . Furthermore, we can compare the network degree of Acs in the GCN and PCN. The network degree of a protein in the GCN (denoted as  $k_{\text{GCN}}$ ) is the number of taxa owning the capacity to express the protein, while the network degree of a protein in the PCN ( $k_{\text{PCN}}$ ) is the number of taxa that have truly expressed the protein. Here,  $k_{\text{GCN}} = 2$  and  $k_{\text{PCN}} = 1$ . Of course, not every protein is ecologically selected. For example, proteins involved in the aminoacyl-tRNA synthesis, critical for the growth of microbes, are not ecologically selected because the loss of ability to synthesize aminoacyl-tRNA inevitably stops the growth of microbes<sup>15,16</sup>. Hence, for each of the proteins involved in aminoacyl-tRNA synthesis, we expect  $k_{\text{GCN}} = k_{\text{PCN}}$  and  $\text{FR}_g = \text{FR}_p$ .

## Illustration of our computational pipeline using a hypothetical community

To illustrate our computational pipeline, let's consider a simple hypothetical example with two species (pink oval vs yellow indented oval in Fig. 1a, b). For the pink species to grow, it can either use the red resource (red pentagon in Fig. 1a) or the blue resource (blue triangle in Fig. 1a) and convert either of them to the green metabolite (green circle in Fig. 1a), which can then be assimilated into the cell biomass. For the yellow species, its growth will only occur by transforming the red resource into the green one to fuel the biomass synthesis (Fig. 1a). If the two species are co-cultured in the same environment to compete for externally supplied red and blue resources, an ideal scenario for them to coexist is that the pink species would choose to consume the blue resource, preventing resource competition with the yellow species (Fig. 1b), similar to the niche partitioning observed in the community of two coexisting *E. coli* strains: CV101 and CV103<sup>21,22</sup>.

We can capture this hypothetical scenario of selective expression mathematically using the GCN and PCN of this community. The bipartite graph and incidence matrix representations of the GCN (or PCN) are shown in Fig. 1a (or Fig. 1b), respectively. Simply comparing the structure of the GCN and the PCN already offers us some insights into ecological niches and metabolic essentiality. For example, let's consider the protein responsible for converting red resource to green metabolite (this protein is represented as the red broken circle in Fig. 1a, b), its degree in the GCN is  $k_{GCN} = 2$ , while its degree in the PCN is  $k_{PCN} = 1$ . This degree reduction is due to distinct ecological niches being occupied by two species when they are cocultured. By contrast, the protein responsible for the assimilation of critical green metabolites (green broken circle in Fig. 1a, b) into biomass does not show a degree reduction from the GCN to the PCN, because it is essential for microbial growth.

An ecologically meaningful approach to understanding the selective expression of different proteins would be to systematically compare their respective  $k_{GCN}$  and  $k_{PCN}$  (Fig. 1c), which are independent of microbial compositions; or their respective  $FR_g$  and  $FR_p$  (Fig. 1d), which naturally involve microbial compositions in the calculation. Consider three distinct protein function types: (1) "niche functions" that are under strong ecological competition (e.g., red broken circle in Fig. 1c, d); (2) "specialist functions" that are specialized by a few taxa (e.g., blue broken circle in Fig. 1c, d); and (3) "essential functions" that are metabolically indispensable for many taxa (e.g., green broken circle in Fig. 1c, d). We anticipate that the three function types will occupy different regions in the  $k_{GCN}$  vs.  $k_{PCN}$  plot (or the  $FR_g$  vs.  $FR_p$  plot). Specifically, for essential functions, both their  $k_{GCN}$  and  $k_{PCN}$  (or  $FR_g$  and  $FR_p$ ) are high.

For specialist functions, both their  $k_{GCN}$  and  $k_{PCN}$  (or  $FR_g$  and  $FR_p$ ) are low. Niche functions have high  $k_{GCN}$  but low  $k_{PCN}$  (or high  $FR_g$  but low  $FR_p$ ).

## Validate our computational pipeline using a consumer-resource model

Note that previously developed Consumer-Resource models (CRMs) only focus on physiologies of microbes (i.e. phenotypes)<sup>29–31</sup>. Simply put, those models ignored genomic capacity or potential functions, but only considered expressed functions (e.g., how species consume different resources). There was no attempt of building a consumer-resource model of microbial communities that integrates both potential and expressed functions. As a first step toward this direction, we constructed such a model.

We assumed three types of protein functions: niche functions (colored red), specialist functions (colored blue), and essential functions (colored green) in a functional pool. For simplicity, each of the niche (or specialist) functions is modeled as the consumption of a unique and externally supplied resource (Fig. 2a1). To model the difference between niche and specialist functions, we assume they are associated with different numbers of species (i.e., “consumers” in the consumer-resource modeling framework). The former should be associated with much more species than the latter. The loss of a niche or specialist function would make a species unable to consume the corresponding externally supplied resource (Fig. 2a2, a3). The loss of an essential function is simply modeled as the reduction of a species’ growth rate (Fig. 2a4). Mathematically, we multiply the intrinsic growth rate of a species by a diminishing factor  $\gamma = 0.95$  for each missing essential function.

The key issue in this genome-aware consumer-resource modelling framework is to decide how microbes select a subset of their potential functions to express. To tackle this issue, we first assigned potential functions to each species (Fig. 2b, left). In particular, for each species, each niche (specialist, or essential) function was assigned to the species’ genome with probability  $p_n$  ( $p_s$ , or  $p_e$ ), respectively. In our simulations, we set  $p_n = p_e = 0.7$  to ensure that we cannot distinguish niche functions from essential functions only based on GCN and thus would like to see if they show different patterns after the community assembly. We set  $p_s = 0.2 < p_n = p_e$  so that specialist functions were assigned to fewer species than niche and essential functions. Then for each species, we determined its truly expressed functions by randomly sub-sampling a subset of its potential functions (Fig. 2b, middle). For function type- $\alpha$  ( $\alpha = 1,2,3$ ), this was achieved by expressing each potential function with a species-specific and function-type-specific probability  $p_{i,\alpha}$  randomly drawn from a uniform distribution  $\mathcal{U}(0,1)$ .



Since different species have different sub-sampling probabilities, some species will tend to be generalists (or specialists). Similar to all consumer-resource models<sup>29–31</sup>, we assume a fixed expression pattern for each species and all resources being supplied so that we don't have to consider the complexity of adaptive expression (such as different expression patterns when different resources are supplied). In the end, we assembled all species in the same community and ran consumer-resource dynamics until the system reached a steady state, for which we constructed the PCN of the survived species (Fig. 2b, right).

We assumed the species pool consists of  $N = 10,000$  species, and the function pool consists of 20 functions for each of the three function types. We introduced 10,000 species to ensure the number of initial species in the assembly simulation is much larger than the number of functions so that we can assemble a high-diversity community in the end. The GCN of the initial species pool is shown in Fig. 2c (left). For each species, we randomly sub-sampled a subset of potential functions to express (middle panel, Fig. 2c). For each species, its true consumption rates are its maximal consumption rates divided by the number of resources the species can use (see Methods) to prevent the selection of generalist species that consume all resources without a penalty<sup>32,33</sup>. Due to the competitive exclusion principle<sup>34</sup>, the maximal number of species survived in the final steady state is 40, because there are 40 unique externally supplied resources ("nutrients") in our model.

In Fig.2c (right), we show a simulation example with 35 species survived in the final steady state. For this assembled steady-state microbial community, we found that the three modeled protein functions types were correctly revealed as three clusters by the Gaussian mixture model in both the comparison of network degree (Fig. 2d) and FR (Fig. 2e). In particular, for niche functions (red cluster in Fig. 2d, e), their mean degree in PCN (2.1) is much lower than that in GCN (24.45), and their mean  $FR_p$  (0.005) is also much lower than their mean  $FR_g$  (0.48). For essential functions (green cluster in Fig. 2d, e), their mean degree in PCN (23.7) is close to that in GCN (26.7), and their mean  $FR_p$  (0.47) is also similar to their mean  $FR_g$  (0.57). For specialist functions (blue cluster in Fig. 2d, e), both their  $k_{GCN}$  and  $k_{PCN}$  (or  $FR_g$  and  $FR_p$ ) are low.

The three functional clusters revealed by the classification of network degrees and functional redundancies for all modeled protein functions exactly match the three function types in our model. Moreover, the relative positioning of the three functional clusters based on our simulation data agrees well with our hypothesis shown in Fig. 1. This clearly validates our hypothesis that niche-occupying proteins have a larger difference in FR and network degree than metabolically essential proteins.

We emphasize that the three functional clusters observed in the  $k_{GCN}$  vs.  $k_{PCN}$  (or the  $FR_g$  vs.  $FR_p$ ) plot is highly nontrivial. It is a result of the community assembly. To demonstrate the importance of community assembly, we randomly picked 35 species (same as the number of survived species) from the initial pool with equal abundances (i.e., the relative abundance is 1/35 for each species) without natural selection and found that it is impossible to distinguish niche functions from essential functions (Fig. 2f, g). Interestingly, for essential functions, we noticed that those species survived after the community assembly tend to have much larger  $FR_p$  (with mean 0.478) than randomly selected species (with mean 0.132). By contrast, for niche functions, survived species tend to have a smaller  $FR_p$  (with mean 0.005) than randomly selected species (with mean 0.133). Similarly, we also computed FR for the same randomly picked 35 species that share the abundances as survived species in the simulation. Again, we cannot differentiate niche functions from essential functions (Supplementary Fig. 1).

We also simulated another community with 100 niche functions, 100 specialist functions, and 100 essential functions. The species pool still consists of  $N = 10,000$  species. As shown in Supplementary Fig. 2), the results are similar to that for the community with fewer functions (Fig. 2).

### Three protein functional clusters observed in human gut microbiomes

After the validation of our computational pipeline using simulated data, we further validated it on real data of human mucosal-luminal interface samples collected from the ascending colon of four children<sup>14,28</sup>. Here we focused on the genus level and annotated the identified proteins from metagenomics and metaproteomics data via the COGs (Clusters of Orthologous genes) database<sup>35,36</sup>. We constructed the GCN and PCN for all the samples following the same procedure as reported in a previous study<sup>14</sup>, and took the intersected COGs between the two networks. In the main text, we focus on the analysis and discussion of subject HM454, and similar findings from the other three subjects are shown in Supplementary Figs. 4-6. For HM454, we used MetaPhlAn2<sup>37</sup> to obtain the taxonomic profile, which includes 85 genera with assigned relative abundances. Raw metagenomic reads and unique peptide sequences detected in metaproteomics were searched against an integrated gene catalog (IGC) database of the human gut microbiome<sup>38</sup> to generate the GCN and PCN respectively. Taxonomic assignment was performed using the 'protein-peptide bridge' method as described previously<sup>14</sup>. More details about data processing can be found in Methods. And the number of intersected COGs for the GCN and PCN associated with HM454 is 1,542. The genus- and COG-level GCN



and PCN of this microbiome sample are shown in Fig. 3a, b. The connectance (i.e., the number of edges divided by the maximal number of possible edges) of the GCN (or PCN) is 0.220 (or 0.049), respectively. The GCN is nested with the nestedness value of 0.667 based on the classical NODF (Nestedness based on Overlap and Decreasing Fill) measure<sup>39</sup> (Fig. 3a; see Methods for details). The PCN has a lower nestedness value of 0.453 for the NODF measure (Fig. 3b).

By comparing the network degree and functional redundancy of one COG in the GCN (one column in Fig. 3a) with those for the same COG in the PCN, we can look into how the COG impacts and is influenced by their metabolic essentiality and connection to occupy ecological niches. For example, COG0539 is the ribosomal protein S1, which has been shown to be essential for some microbes<sup>40–44</sup>. For subject HM454, 20 genera have COG0539 in the GCN, while 15 genera have this COG in the PCN, hence  $k_{GCN} = 20$  and  $k_{PCN} = 15$ . Additionally, COG0539 has a similar level of functional redundancy in GCN and PCN:  $FR_g = 0.476$  and  $FR_p = 0.461$ . These results suggest that COG0539 is crucial for microbial metabolism, and not ecologically selected. Another example that falls into a different category (i.e., niche functions) is COG1116, which is the ABC-type nitrate/sulfonate/bicarbonate transport system<sup>45</sup>. For COG1116, we have  $k_{GCN} = 22 \gg k_{PCN} = 2$ ; and  $FR_g = 0.388 \gg FR_p = 0.004$ , which is evidence for the further specification in transporting nitrate, sulfonate, or bicarbonate across community members on the protein level. Different from the previous examples, some functions are specialized by a few genera on the gene level and thus are still specialized by those genera on the protein level. For example, COG1018 (Ferredoxin-NADP reductase), which has  $k_{GCN} = k_{PCN} = 1$  and  $FR_g = FR_p = 0.0$ , is classified as a specialist function.

To systematically explore the difference between GCN and PCN, we visualized the difference in the network degree (Fig. 3c) and functional redundancy (Fig. 3d) for all COGs. As can be seen in Fig. 3c for comparing network degrees, nearly all COGs are below the black dashed line of  $k_{GCN} = k_{PCN}$  because the map from the genomic capacity to protein function is a sub-sampling process. The network degrees in PCN for almost all points are less than 10 (1,365 out of 1,542) and much less than their corresponding network degrees in GCN (349 out of 1,542 COGs have network degrees less than 10). 804 of 1,542 COGs have a reduction in network degree by more than 80%. Eventually, the major difference in network degree will lead to a significant difference in functional redundancy, although the reduction in network degree from GCN to PCN cannot fully explain why many COGs have  $FR_p \sim 0$  (744 out of 1,542 have

FR<sub>p</sub> < 0.01 in Fig. 3d). Indeed, the network degrees for COGs in the PCN positively correlate with FR<sub>p</sub>, but there is no simple relationship between k<sub>PCN</sub> and FR<sub>p</sub> (Fig. 3e). For example, for L-arabinose isomerase (COG2160), its network degree in GCN (k<sub>GCN</sub> = 8) is fairly close to the network degree in PCN (k<sub>PCN</sub> = 7), but its FR<sub>p</sub> (0.04) is much lower than FR<sub>g</sub> (0.23) since the genus *Blautia* (which makes up 22% of the subject HM454's total microbial abundance) didn't express L-arabinose isomerase, even if it has this capacity encoded in its genome.

We applied the Gaussian mixture model fitted on simulated data to classify all protein functions in the real data and obtained 3 clusters from both the k<sub>GCN</sub> vs. k<sub>PCN</sub> plot (Fig. 3c) and the FR<sub>g</sub> vs. FR<sub>p</sub> plot (Fig. 3d). Despite that the clustering of protein functions in real data looks weaker than that in simulated data, the relative positioning of the three clusters (shaded areas in Fig. 3c, d) agree well with our hypothesis shown in Fig. 1, as well our simulation results shown in Fig. 2. We suspect that the weakened clustering might be due to (1) the variation of k<sub>GCN</sub> (or FR<sub>g</sub>) in real data (Fig. 3c, d) is much larger than that in simulated data (Fig. 2d, e); the low resolution of the GCN and PCN in the real data (both were constructed at the genus level).

Note that some points in Fig. 3c, d are above the diagonal line, contradicting the subsampling argument for the gene expression. For instance, we noticed that for the subject HM454, 12 genera have COG0094 in the GCN, while 25 genera have this COG in the PCN. Additionally, COG0094 is even less redundant in the GCN (FR<sub>g</sub> = 0.166) than it is in the PCN (FR<sub>p</sub> = 0.641). FR<sub>g</sub> should be always larger than FR<sub>p</sub> if the PCN was a proper subgraph of the GCN for COG0094. We believe this contradiction is largely due to the metagenomic sequencing depth and the metaproteomic identification depth. We know that both metagenomics and metaproteomics have depth limitations and require sufficient depth to detect genes or proteins, respectively. More specifically, some proteins detected by the ultra-deep metaproteomics are not found in putative protein sequences annotated from metagenomes. For example, if more proteins were assigned to one COG by the metaproteomics than annotated metagenomes, it indicates the number of taxa that express proteins belonging to the COG is higher than the number of taxa that own the COG. As a result, the network degree of the COG in the GCN is even higher than its network degree in the PCN, making FR<sub>p</sub> of the COG larger than its FR<sub>g</sub> (evidenced by COG0094).

### Comparing FR<sub>g</sub> with FR<sub>p</sub> pinpoints ecological niches and metabolic essentiality

In order to justify whether or not the FR comparison for many COGs is ecologically or metabolically meaningful, we focus on two types of proteins: ABC-type transporters (under

strong ecological selection because they directly influence the ecological interactions and are influenced by resource availability)<sup>45–47</sup> and ribosomal proteins (under weak ecological selection because of their essentiality)<sup>42–44</sup>.

ABC-type transporters are energy-requiring transporter proteins responsible for obtaining and releasing resources in the environment<sup>45–47</sup>. For example, if we consider a particular transporter responsible for the uptake of glucose from the environment, theoretically only top consumers of glucose would have the chance to claim this niche (consumption of glucose) from the ecological standpoint. Consequently, we should expect a specification in glucose consumption on the level of protein functions, even though many species have the capacity to utilize it. For the gut microbiota sample we investigated, we indeed found that  $k_{GCN}$  for all ABC-type transporters are much larger than their  $k_{PCN}$  (Fig. 4a). Similarly, we also found that  $FR_g$  for all ABC-type transporters are much larger than their  $FR_p$  (Fig. 4b). Many transporter proteins were classified to the red cluster (i.e., the cluster of niche functions) in Fig. 4b. Some transporter proteins were classified to the blue cluster (i.e., the cluster for specialist functions) due to the specialization on the gene level. As a result, such specialization would be carried to the protein level. Some transporter proteins were classified to the green cluster (i.e., the cluster for essential functions) because they have been proven essential for microbes. One example is the ABC-type Fe<sup>3+</sup>/spermidine/putrescine transporter (COG3842) which has  $FR_g = 0.339$  and  $FR_p = 0.285$ . It has been shown that iron is essential for bacteria as it functions as a co-factor in iron-containing proteins in redox reactions, metabolic pathways, and electron transport chain mechanisms<sup>48,49</sup>.

Ribosomal proteins are necessary for the growth of all living organisms because, as we know, the ribosome is the place where other proteins are synthesized<sup>50,51</sup>. Since ribosomal proteins are an indispensable part of microbial survival, all abilities of synthesizing such proteins are expected to be expressed. In our data, many ribosomal proteins were classified to the green cluster (i.e. the cluster for essential functions). Moreover, we found that their  $k_{GCN}$  were very close to their  $k_{PCN}$  (Fig. 4e). In Fig. 4f, we compared  $FR_g$  with  $FR_p$  and found many ribosomal proteins were classified to the green cluster (i.e. the cluster for essential functions), agreeing with our expectation that proteins with high  $FR_g$  and  $FR_p$  are more likely to be essential functions. Interestingly, two ribosomal proteins (L28 and L34) colored red in Fig. 4e have been shown to be non-essential<sup>41,42,52</sup> to microbes such as *E. coli*. Some specialized ribosomal proteins in microbial genomes continue to be specialized on the protein level and thus were classified to the blue cluster (i.e., the cluster for specialist functions).

Alternatively, we looked at the distribution of network degrees (Fig. 4c, g) and the distribution of functional redundancy ( $FR_g$  or  $FR_p$  in Fig. 4d, h) for the two protein types to observe their difference. For ABC-type transporters, the distribution of network degrees in PCN is close to 0 (having a median of 2), while the median of network degrees in GCN is 25. For ribosomal proteins, the distribution of network degrees in PCN (median is 12) is similar to that in GCN (median is 14). For ABC-type transporters, the distribution of  $FR_p$  in PCN is close to 0 (with a median  $\sim 0.01$ ), while the median of  $FR_g$  in GCN is around 0.30. For ribosomal proteins, the distribution of  $FR_p$  in PCN (median  $\sim 0.20$ ) is similar to the distribution of  $FR_g$  in GCN (median  $\sim 0.21$ ). The same patterns of ABC transporters showing a big reduction (in functional redundancy and network degree) and ribosomal proteins showing little difference are also true for the other 3 individuals (Supplementary Figs. 9-11).

We also validated the above results using a different functional annotation method, KEGG Orthology (KO)<sup>53-56</sup>. The annotation rate of proteins involved in PCN of the four individual microbiomes is 78% (much lower than 92% which we had for the COG annotation). The contrasting difference between ABC-type transporters and ribosomal proteins is well preserved (see Supplementary Fig. 7). Additionally, the distribution of  $FR_p$  shows a dramatic difference across KO groups (Supplementary Fig. 8). Some ecologically strongly selected KO groups have small  $FR_p$ , while other metabolically essential KO groups show fairly large  $FR_p$  and big variations (see Supplementary Fig. 8). For example, almost all proteins in ABC transporters and PTS systems have  $FR_p$  close to zero (Fig. Supplementary Fig. 8), and transporters and PTS systems are well-known as the ecologically selected groups<sup>45-47,57</sup>. As a comparison, proteins from Aminoacyl-tRNA biosynthesis, glycolysis, and ribosomes all have big  $FR_p$  and huge variations across different proteins within the group (Supplementary Fig. 8). In the past, the metabolic essentiality has been demonstrated for Aminoacyl-tRNA biosynthesis<sup>15,16</sup>, glycolysis<sup>58,59</sup>, and ribosomes<sup>42-44</sup>.

## **The response of community and protein expression to the introduction of sugars**

In ecology, a niche is often defined as an abiotic and biotic factor that supports the survival of species<sup>9,60-62</sup>. Therefore, niche functions are associated with corresponding limiting resources involved in those functions. For example, COG1879 (ABC-type sugar transport system, periplasmic component, contains N-terminal xre family HTH domain) which is categorized as a niche function owing to its high  $FR_g$  of 0.486 and low  $FR_p$  of 0.041 for the subject HM454, is

associated with widely competed sugars in microbial communities. After inferring niche functions such as ABC-type transporters by our computational pipeline, we wonder if it is possible to influence the community structure by externally supplying more limiting resources involved in the niche functions. To demonstrate this, we resort to the *in vitro* community and are interested in how the community structure and expression of proteins involved in niche functions respond to supplied limiting sugars. Specifically, we would like to see how proteins relevant to ecological niche functions within one taxon change their expressions to achieve a better living strategy for the taxon.

We used rapid assay for individual microbiome (RapidAIM)<sup>63</sup>, which maintains the functional profiles of individual gut microbiomes *in vitro*<sup>64</sup>, to culture three different individual human gut microbiota samples, and used metaproteomics to observe how taxon-specific expression of proteins in the niche functional cluster respond to the presence of glucose, fructose and kestose (Fig. 5a). Samples were cultured in technical triplicates, and were taken at 0, 1, 5, 12, and 24 hours of culturing for optical density and metaproteomic analyses. 11-plex tandem mass tag (TMT11plex) was used for metaproteomic quantification<sup>65</sup> for a total of 189 samples. To reflect the effect of introduced sugars on protein expression levels, we used log2 of the ratio between normalized protein abundances/intensities (see Methods for details) in the treatment and that in the control group (i.e. log2 of fold change in Fig. 5). We hypothesized that the excessive supply of sugars renders carbon resources no longer limited and instead microbes start to compete for other resources in relatively short supplies compared to carbon resources such as nitrogen resources or amino acids because microbes need all those resources proportionally (Fig. 5a). Therefore, microbes might have to over-express proteins to uptake more non-carbon limiting resources to achieve better growth.

To understand how each taxon interacts with the environment and how introduced sugars modulate the interaction, we focused on log2 fold changes of ABC-type transporters 5 hours later whose expression levels reveal rates for transporting nutrients (Fig. 5b-d). When glucose is supplied in an excessive amount, log2 fold changes of most COGs are close to zeros except for COG1126 (ABC-type polar amino acid transport system, ATPase component), COG1653 (ABC-type glycerol-3-phosphate transport system, periplasmic component), COG1879 (ABC-type sugar transport system, periplasmic component, contains N-terminal xre family HTH domain), and COG4166 (ABC-type oligopeptide transport system, periplasmic component). Many pronounced changes happen to the genus *Holdemanella* and it is interesting to note that *Holdemanella* reduces the expression of transporters for importing sugars (COG1879) and an intermediate in the glycolysis glycerol-3-phosphate (COG1653)

when glucose is added. Instead, it increases the expression of COG1126 which transports polar amino acids. This strategy benefits *Holdemanella* because the fraction of proteins from *Holdemanella* over all proteins in the community increases from 13.5% in the control to 15.8% with the added glucose. We also measured log2 fold changes of ABC-type transporters when fructose, glucose and fructose, or kestose is added and their overall patterns (Fig. 5c-e) are similar to the pattern when glucose is added (Fig. 5b). The correlation in log2 fold changes of ABC-type transporters between different added sugars is significant (Supplementary Fig. 12). Similar fold changes of ABC-type transporters were observed for metaproteomic measurements 12 hours, and 24 hours later (Supplemental Figs. 14-15), while the fold changes 1 hour later are still fairly small (Supplemental Fig. 13). We also attempted to look at how ribosomal proteins respond to sugar supplies (Supplementary Fig. 16). Overall, log2 fold changes of ribosomal proteins are overwhelmingly positive, which probably implies a faster growth for microbes when simple sugars are supplied<sup>32,33</sup>. Therefore, we demonstrated that the sugars associated with the niche function (i.e., the sugar transport system) can be introduced to influence gene expression and modulate the community structure.

## Discussion

Understanding the functions of proteins in the metabolism and how they are influenced by various ecological interactions is important to fully characterize ecological niches in a given microbial community. Typically, to check if a protein is metabolically essential, one has to knock out the gene in one microbial species that codes for the protein to check how the growth rate of the species reduces<sup>42-44</sup>. A usual way to determine a limiting resource often that is utilized by a protein follows: modify resource supplies and see how the total biomass changes<sup>66-69</sup>. Here, to complement those traditional experimental methods, we proposed a simpler computational method that can identify metabolic and ecological functions of proteins via the comparison of their  $FR_g$  and  $FR_p$ , as well as their  $k_{GCN}$  and  $k_{PCN}$ . We validated this computational method using both model-generated synthetic data and real data for human gut microbiomes. Also, when we selected two types of proteins (ABC-type transporters and ribosomal proteins in the real data representing niche functions and essential functions, respectively), most predicted protein functional clusters of the two types of proteins fell into the niche function cluster and the essential function cluster, respectively. Besides these two protein types, we were able to generate a list of  $FR_p$  and  $FR_g$  for all COGs (see Supplemental Data 1-4), which is useful for understanding the metabolic and ecological functions of proteins.



The presented results help us reconcile the conflict between the niche theory in ecology<sup>62,70,71</sup> and the observed functional redundancy<sup>11,12</sup>. The traditional niche theory is grounded in the competitive exclusion principle, stating that a resource (or niche) cannot be occupied by two species (or more than two species) for the steady-state conditions<sup>62,70,71</sup>. As a result of competition, organisms within the same community develop different surviving strategies to minimize their competition. One interesting example is the repetitive established coexistence between two evolved *E. coli* strains, even though a single clone of *E. coli* is initiated and maintained in a glucose-limited continuous or serial culture<sup>21,72,73</sup>. Cross-feeding between two evolved *E. coli* strains can be established when one bacterial strain consumes overflow metabolites like acetate excreted by the other bacterial strain<sup>21</sup>. Hence, the two strains avoid competition by specification on different resources (glucose and acetate). However, the picture from the niche theory clashes with the observed functional redundancy in microbial communities because the functional redundancy implies that many species own the same functions in their genomes<sup>11,12</sup>. We solved this dilemma by pointing out that proteins related to occupying ecological niches usually have very low  $FR_p$  and large  $FR_g$ . Therefore, if we apply this concept in reverse, then large  $FR_g$  and small  $FR_p$  could help us to pinpoint niche functions.

There is a long-standing gap between the ecological model which considers the protein functions of organisms and the data analysis of genomic data to give ecological insights. Ever since Robert MacArthur proposed a community model in 1970 to consider how different consumers compete exclusively for renewing resources<sup>74</sup>, many extensions of this model were proposed to include more complex ecological factors such as cross-feeding interactions<sup>75–78</sup> and multiple essential nutrients<sup>79</sup>. Almost all of them focus on the phenotype of microbes because only functions of expressed proteins are relevant for the consumption and production of nutrients in the ecosystem. Due to the lack of metaproteomic data, many computational approaches attempting to generate ecological implications rely on the over-complete inferred protein capacity derived from genomes<sup>7,9–12</sup>. To reconcile this gap, we built an ecological framework with the genomic capacity and protein functions together by introducing species with sub-sampled functions. The model framework is useful for explaining the difference between genomic capacity and protein functions. The selective expression can be considered as the same microbe with different expressions under different environments<sup>80–82</sup> or evolved strains from the same species that have distinct metabolic niches observed in evolutionary experiments of microbes<sup>83,21,22</sup>.

It is worth noting that the assumption of the trade-off between generalists and specialists (represented by assuming that the total proteome is relatively constant) is very

important. In our model, this assumption is achieved by considering true consumption rates in PCN as maximal consumption rates in GCN divided by the number of resources. The importance of this trade-off lies in the fact that it forces the niche partitioning among species. In the absence of this assumption, there is no pattern of redundancy difference since generalists can always out-compete specialists. This trade-off makes sense because typically the total proteome budgets for microbes have been observed to be relatively fixed<sup>32,33</sup>.

## Methods

**In-vitro culture of single gut bacterial strains with added sugars.** Five gut commensal bacterial strains, *Bacteroides vulgatus* ATCC 8482, *Bacteroides ovatus* ATCC 8483, *Bacteroides uniformis* ATCC 8492, *Blautia hydrogenotrophica* DSM 10507, *Escherichia coli* DSM 101114 were cultured with or without added sugars (glucose, sucrose and kestose). The base culture medium without sugar added were modified based on the Yeast Casitone Fatty Acids (YCFA) broth, containing 10.0 g/L casitone, 2.5 g/L yeast extract, 45 mg/L MgSO<sub>4</sub>·7H<sub>2</sub>O, 90 mg/L CaCl<sub>2</sub>·2H<sub>2</sub>O, 450 mg/L K<sub>2</sub>HPO<sub>4</sub>, 450 mg/L KH<sub>2</sub>PO<sub>4</sub>, 900 mg/L NaCl, 1.0 mg/L resazurin, 4.0 g/L NaHCO<sub>3</sub>, 1.0 g/L L-Cysteine-HCl, 10 mg/L Hemin, 1.90 mL/L acetic acid, 0.7 mL/L propionic acid, 90 µL/L iso-butyric acid, 100 µL/L n-valeric acid, 100 µL/L iso-valeric acid, 0.02 mg/L biotin, 0.02 mg/L folic acid, 0.05 mg/L thiamine-HCl, 0.05 mg/L riboflavin, 0.001 mg/L vitamin B12, 0.05 mg/L aminobenzoic acid. The pH was adjusted to between 6.7-6.8, and autoclaved media were pre-reduced in an anaerobic chamber overnight. 5 g/L of different sugars (glucose, sucrose, and kestose) were added to the base medium as treatment groups. Master tubes of single bacterial strains were first cultured on Tryptic Soya Agar (TSA) containing 5% sheep blood using the streak plate method. A single colony was picked from each agar plate and inoculated into the base culture medium to culture for 24 hours, before inoculating 100 µL of each culture into 10 mL of four different media: base medium without sugar added, with glucose added, with sucrose added and with kestose added. After culturing for 24 hours, optical density at 600 nm was tested in technical triplicates for each sample. Cultured microbial cells were purified by washing with phosphate buffered saline (PBS) buffer three times, and the resulting microbial pellets were stored at -80 °C for proteomics analysis.

**In-vitro human gut microbiota culture with added sugars.** Three healthy individual microbiota samples were collected and biobanked using our live microbiota biobanking protocol<sup>84</sup>. The study was approved by the Ottawa Health Science Network Research Ethics Board at the Ottawa Hospital, Ottawa, Canada (# 20160585–01 H). The frozen microbiome samples were thawed at 37 °C and cultured in our optimized culture medium<sup>64</sup> with or without the presence of different sugars (10 mM glucose, 20 mM fructose, 10 mM glucose + 20 mM fructose, or 10 mM kestose). Samples were cultured in technical triplicates, and were taken at 0 hr, 1hr, 5 hr, 12 hr and 24 hr of culturing for optical density and metaproteomic analyses. After culturing, 96-well deep well plates were first centrifuged at 3,000 g for 45 min under 4 °C. Then the pellets were washed in 4 °C phosphate buffered saline (PBS) buffer and centrifuged at 3,000 g for 45 min again, before pelleting and removing culture debris three times using 300 g, 4 °C, 5 min centrifugation. Microbial suspensions were then centrifuged at 3,000 g, 4 °C for another 45 min. The purified cell pellets were stored at -80 °C before protein extraction.

**Protein extraction, digestion and LC-MS/MS analysis.** For single strain samples, proteins were extracted with 4% SDS 8M urea buffer in 100 mM Tris-HCl buffer and precipitated overnight at -20 °C, before being purified by washing with ice-cold acetone three times. Quantified proteins were then reduced and alkylated before being digested using trypsin (50:1 protein-to-trypsin ratio) for 24 hours at 37 °C and were desalted using reverse phase beads<sup>85</sup>. Proteomic samples were analyzed using an Orbitrap Exploris 480 mass spectrometer (ThermoFisher Scientific Inc.) coupled with an UltiMate 3000 RSLCnano liquid chromatography system following a 1-hour gradient of 5 to 35% (v/v) acetonitrile (v/v) at the flow rate of 300 L/min. MS full scan was performed from 350 - 1400 m/z with a resolution of 60,000, followed by an MS/MS scan of 12 most intense ions, a dynamic exclusion repeat count of one, exclusion duration of 30 s, and resolution of 15,000. Metaproteomics samples of the cultured individual microbiomes were prepared using a semi-automated approach. Briefly, samples were lysed in a buffer containing 8 M urea, 4% SDS in 100 mM Tris-HCl (pH = 8.0) to extract microbial total proteins. The proteins were purified by a double-precipitation procedure in 50%:50%:0.1% (v/v/v) acetone: ethanol: acetic acid solution. Protein digestion and desalting steps were performed using an automated liquid handler (Hamilton Nimbus-96). Briefly, 100 µg proteins were dissolved in 100 µL 6 M urea in 100 mM Tris-HCl (pH 8) buffer, before being reduced by 10 µL 0.1 M dithiothreitol (DTT) solution under 56 °C for 30 minutes and alkylated by 10 µL 0.2 M iodoacetamide (IAA) solution in dark, 25 °C for 40 minutes. Samples were each added 1000 µL 100 mM Tris-HCl buffer containing 2 µg/mL trypsin (trypsin:proteins = 1:50) for

a 24-hour digestion under 37 °C, before being desalted using an automated pipeline based on reverse-phase (RP) desalting columns. 11-plex tandem mass tag (TMT11plex) was used for metaproteomic quantification for a total of 189 samples. An even mixture of all samples was used as the reference channel in each 11-plex. Samples were scrambled before labeling with TMT11plex, so that each labeled sample contains samples from different individuals, different time points and different treatments to avoid any bias that may be induced between analyses. TMT-labelled samples were analyzed using an Orbitrap Exploris 480 mass spectrometer (ThermoFisher Scientific Inc.) coupled with an UltiMate 3000 RSLCnano liquid chromatography system following a 2-hour gradient of 5% to 35% solvent B (80% acetone nitrile, 0.1% formic acid, v/v).

**Datasets.** Metagenomics data corresponding to the ultra-deep metaproteomic analysis of the four individual microbiomes were obtained from the previous MetaPro-IQ study<sup>14,28</sup> (accessible from the NCBI sequence read archive (SRA) under the accession of SRP068619) and the same samples were reanalyzed by an ultra-deep metaproteomics approach<sup>14</sup> (accessible through the ProteomeXchange Consortium (<http://www.proteomexchange.org>) via the PRIDE partner repository<sup>86</sup>). Proteomics dataset of the cultured singles strain samples has been deposited to ProteomeXchange Consortium via the PRIDE partner repository. Metaproteomic dataset of the RapidAIM-cultured microbiome samples has been deposited to ProteomeXchange Consortium via the PRIDE partner repository.

**Database search and data processing.** Proteomics database searches were performed by combining FASTA databases of the individual strains downloaded from NCBI. The databases were combined for performing database search using MaxQuant<sup>87</sup> 1.6.17.0, with the label-free quantification option turned off. Metaproteomic database searches of cultured microbiome samples were performed using MetaLab V2.2<sup>88</sup>, MaxQuant option was used to search the TMT dataset against the IGC database of the human gut microbiome. The resulting data table was normalized using R package MSstatsTMT<sup>89</sup>, and missing values were imputed using R package DreamAI<sup>90</sup>. The "fraction" of each taxon-specific protein is computed by dividing the protein intensity by the sum of intensities of all proteins assigned to the same taxon. The log2 fold change of each protein is obtained by taking log2 of the ratio between its fraction in the treatment group (with added sugars) and its fraction in the control group (without added sugars).

**Generation of GCN and PCN.** For the ultra-deep metaproteomic dataset, the genus-COG version of GCN and PCN tables were directly obtained from the previous work<sup>14</sup>. In addition, here we generated a genus-KEGG version of GCN and PCN for each individual microbiome using a similar method. Briefly, for the genus-KEGG GCN, by searching raw metagenomic reads against an integrated gene catalog (IGC) database of the human gut microbiome<sup>38</sup>, we obtained a list of proteins quantified by read counts. FASTA sequences of these proteins were searched against the KEGG database using GhostKOALA<sup>91</sup>. Taxonomic origination of the proteins was obtained by searching against an in-house database generated with the NCBI non-redundant (nr) database (downloaded 2/3/2016). To generate genus-KEGG PCN, the taxonomic table of the metaproteomics dataset was directly obtained from MetaLab, and KEGG annotation was also performed by querying protein FASTA sequences with GhostKOALA. Protein group intensity was used as the quantification information in PCNs. For the proteomic dataset of single strains, the whole proteomic FASTA database was submitted to EggNOG mapper (<http://eggno-mapper.embl.de/>, submitted Oct-30-2021, ran emapper.py 2.1.6) to obtain functional annotations. To generate GCN, protein coding sequence (CDS) files were downloaded from NCBI, and the count of each protein id in the CDS files was considered as the copy number of each gene in the GCN. For PCN generation, intensities of identified proteins matched to each strain were used. Note that protein ids in the CDS file were 100% matched with those in the proteomic FASTA database in each strain. For the metaproteomics dataset of the cultured microbiome samples, functional information for the generation of PCN was obtained from the resulting functional table automatically generated by the MetaLab software. Taxonomic assignment was performed using the ‘protein-peptide bridge’ method as described previously<sup>14</sup>. The PCNs for this dataset were then generated based on intensities of COG-genus pairs.

**Normalized gene-level functional redundancy (nFR<sub>g</sub>) and normalized protein-level functional redundancy (nFR<sub>p</sub>).** Across multiple samples, it is pointless to compare the FR<sub>g</sub> or FR<sub>p</sub> directly because of the difference in microbial taxonomic diversities. In fact, it has been shown in the past that the normalized functional redundancy, which is the functional redundancy divided by the taxonomic diversity, can be compared across samples<sup>12</sup>. In our study, the definition for nFR<sub>g</sub> is

$$\text{nFR}_g = \frac{\sum_{i=1}^N \sum_{j \neq i}^N (1 - a_{ij}^{\text{GCN}}) p_i p_j}{\sum_{i=1}^N \sum_{j \neq i}^N p_i p_j}, \quad (1)$$

and the definition for  $nFR_p$  is

$$nFR_p = \frac{\sum_{i=1}^N \sum_{j \neq i}^N (1 - d_{ij}^{PCN}) p_i p_j}{\sum_{i=1}^N \sum_{j \neq i}^N p_i p_j}. \quad (2)$$

## The community assembly model.

*Step 1: Assignment of species' genomic capacity.* Three types of protein functions are modeled: niche function, specialist function, and essential function. Both specialist function and niche function are considered as the capacity to consume a unique and externally supplied resource. The probability of a random consumer being assigned the ability to have a niche function is 0.7. To make fewer species own specialist functions in their genomes, the probability of a random consumer being assigned the ability to have a specialist function is 0.2, much lower than the probability of owning a niche function. The maximal consumption rate of a resource by one species represents the consumption rate that the species would have if it allocates the entire proteome (100%) to the consumption of the resource. If many resources are consumed, the total proteome has to be divided into several parts and the consumption rates would be a fraction of the corresponding maximal consumption rates. The essential function is not modeled as the consumption of alternative resources due to its metabolic essentiality. Instead, the essential function is modeled as multiplying the growth rate by a factor of 0.95 for each missing essential function.

*Step 2: Assignment of species' protein functions based on their genomic capacity.* Each species sub-samples its genomic potential functions with a sub-sampling probability  $p$  (which is a random number uniformly distributed between 0 and 1) to obtain its protein functions (i.e. which resource it can truly consume). As a result, all protein functions of species form the basis for PCN. The true consumption rate of one species on a resource is its maximal consumption rate on the resource divided by the number of resources that can be utilized by the species. This process can be thought of as the proteome allocation to consume several resources simultaneously<sup>32,33</sup>. This assumption imposes a trade-off between a generalist and a specialist species: a generalist species utilizes more resources but has lower consumption rates for all resources, while a specialist species consumes fewer resources but has higher consumption rates for consumed resources.

*Step 3: Community assembly.* We assumed a chemostat environment, similar to the setting considered by many Consumer-Resource models<sup>75,77</sup>. The dilution rate  $D$  is considered as 0.1 per hour. A fixed number of resources is considered and the pool concentrations (or supply rates) for all resources are assumed to be the same for simplicity. For each species, the growth



rate is treated as the sum of consumption rates for different resources divided by the yield. For simplicity, all yields are assumed to be equal ( $Y = 1$ ). Overall, the dynamics for the concentrations of resource  $i$  (denoted as  $C_i$ ) and the abundance of the species  $\alpha$  (written as  $B_\alpha$ ):

$$\frac{dC_i}{dt} = h_i - DC_i - \frac{\sum_{\beta} a_{\beta i} Y^{N_m} B_{\beta} C_i}{Y}, \quad (3)$$

$$\frac{dB_{\alpha}}{dt} = -DB_{\alpha} + \sum_j a_{\alpha j} \gamma^{N_m} B_{\alpha} C_j, \quad (4)$$

where  $a_{\alpha i}$  is the consumption rate of species  $\alpha$  on resource  $i$ ,  $h_i$  is the supply rate of resource  $i$ ,  $Y$  is the same yield assumed for all resources,  $\gamma (= 0.95)$  is the diminishing rate for the overall consumption rate that is multiplied for each missing essential function, and  $N_m$  is the number of missing essential functions. The consumption rate of one species of a resource is randomly drawn from the uniform distribution between 0 and 1. Eventually, for each species, its true consumption rates are its randomly drawn consumption rates divided by the number of resources the species can use to constrain the total proteome budget<sup>32,33</sup>. The incidence matrix of the consumption abilities establishes part of PCN for niche functions and specialist functions of the species. The entire PCN is completed by including the presence/absence information of all essential functions.

*Step 4: Generate GCN and PCN for survived species.* When we simulated the above community assembly process to reach a steady-state in the chemostat environment, survived species can be found as species existing with non-negative abundances at the end of the simulation. For survived species, we can reconstruct the GCN and PCN for them. Within equipped GCN and PCN, we would be able to compute  $FR_g$ ,  $FR_p$ , and network degrees ( $k_{GCN}$  and  $k_{PCN}$ ).

**Calculation of nestedness.** To reveal the nested structure of an incidence matrix, we first need to use the Nestedness Temperature Calculator (NTC)<sup>92</sup> to organize the matrix. Then we adopted the NODF (Nestedness based on Overlap and Decreasing Fill) measure previously defined<sup>39</sup>. The measure can only be computed for binary incidence matrices. As with any perfectly nested matrix, two properties must be present: (1) decreasing fill, which means that the columns below and to the right should have fewer entries than the columns above and to the left; and (2) paired overlap, which implies that when an entry appears in the columns below and to the right, it should also appear in the columns above and to the left. The NODF measure is calculated by averaging these two properties across all pairs of an upper and lower row and a left and right column. For the comparison of each pair, if decreasing fill is not

satisfied, the pair will contribute 0 to the total nestedness. Otherwise, the pair's contribution is the percentage overlap in non-zero entries between the two rows or two columns.

**Statistics.** To calculate correlation throughout the study, we used Pearson's correlation coefficient. Wherever we used *P* values, we explained in the Methods how we calculated them, since for all such measurements in the study, we calculated the associated null distributions from scratch. All statistical tests were performed using standard numerical and scientific computing libraries in the Python programming language (version 3.7.1) and Jupyter Notebook (version 6.1).

**Data and code availability.** All code for simulations used in this manuscript can be found at XXX.

**Acknowledgements.** We thank Janice Mayne for the help with providing biobanking samples and thank Zhibin Ning for the help with running mass spectrometry. D.F. acknowledges grants from Natural Sciences and Engineering Research Council of Canada (NSERC), and the Government of Canada through Genome Canada and the Ontario Genomics Institute (OGI-114 & OGI-149). D.F. acknowledges a Distinguished Research Chair from the University of Ottawa. Y.-Y.L. acknowledges grants from the National Institutes of Health (R01AI141529, R01HD093761, RF1AG067744, UH3OD023268, U19AI095219, and U01HL089856).

**Author contributions.** Y.-Y.L. and D.F. supervised the study. T.W. and Y.-Y.L. conceived the project. All authors designed the research. L. L. prepared and curated the empirical data as well as performed all wet-lab experiments. T.W. analyzed all data and developed the ecological model. T.W. wrote the initial manuscript. All authors edited and approved the manuscript.

**Competing Interests.** The authors declare no competing interests. D.F. co-founded MedBiome Inc., a clinical microbiomics company.

## References

1. Tyson, G. W. *et al.* Community structure and metabolism through reconstruction of microbial genomes from the environment. *Nature* **428**, 37–43 (2004).
2. Flint, H. J., Scott, K. P., Louis, P. & Duncan, S. H. The role of the gut microbiota in nutrition and health. *Nature Reviews Gastroenterology & Hepatology* **9**, 577–589 (2012).
3. Lloyd-Price, J. *et al.* Multi-omics of the gut microbial ecosystem in inflammatory bowel diseases. *Nature* **569**, 655–662 (2019).
4. Paerl, null & Pinckney, null. A Mini-review of Microbial Consortia: Their Roles in Aquatic Production and Biogeochemical Cycling. *Microb Ecol* **31**, 225–247 (1996).
5. Falkowski, P. G., Fenchel, T. & Delong, E. F. The Microbial Engines That Drive Earth's Biogeochemical Cycles. *Science* **320**, 1034–1039 (2008).
6. Louca, S., Parfrey, L. W. & Doebeli, M. Decoupling function and taxonomy in the global ocean microbiome. *Science* **353**, 1272–1277 (2016).
7. Levy, R. & Borenstein, E. Metabolic modeling of species interaction in the human microbiome elucidates community-level assembly rules. *PNAS* **110**, 12804–12809 (2013).
8. Pacheco, A. R., Moel, M. & Segrè, D. Costless metabolic secretions as drivers of interspecies interactions in microbial ecosystems. *Nature Communications* **10**, 103 (2019).
9. Fahimipour, A. K. & Gross, T. Mapping the bacterial metabolic niche space. *Nature Communications* **11**, 4887 (2020).
10. Louca, S. *et al.* High taxonomic variability despite stable functional structure across microbial communities. *Nat Ecol Evol* **1**, 15 (2016).
11. Louca, S. *et al.* Function and functional redundancy in microbial systems. *Nature Ecology & Evolution* **2**, 936 (2018).
12. Tian, L. *et al.* Deciphering functional redundancy in the human microbiome. *Nature Communications* **11**, 6217 (2020).
13. Franzosa, E. A. *et al.* Relating the metatranscriptome and metagenome of the human gut. *Proc Natl Acad Sci U S A* **111**, E2329–2338 (2014).
14. Li, L. *et al.* Revealing Protein-Level Functional Redundancy in the Human Gut Microbiome using Ultra-deep Metaproteomics. *bioRxiv* 2021.07.15.452564 (2021) doi:10.1101/2021.07.15.452564.
15. Ibba, M. & Soll, D. Aminoacyl-tRNA synthesis. *Annu Rev Biochem* **69**, 617–650 (2000).
16. Parker, D. J. *et al.* Growth-Optimized Aminoacyl-tRNA Synthetase Levels Prevent Maximal tRNA Charging. *Cell Syst* **11**, 121–130.e6 (2020).
17. Jacob, F. & Monod, J. Genetic regulatory mechanisms in the synthesis of proteins. *J Mol Biol* **3**, 318–356 (1961).
18. Okano, H., Hermesen, R., Kochanowski, K. & Hwa, T. Regulation underlying hierarchical and simultaneous utilization of carbon substrates by flux sensors in *Escherichia coli*. *Nat Microbiol* **5**, 206–215 (2020).
19. Kumari, S. *et al.* Regulation of Acetyl Coenzyme A Synthetase in *Escherichia coli*. *Journal of Bacteriology* **182**, 4173–4179 (2000).
20. Starai, V. J. & Escalante-Semerena, J. C. Acetyl-coenzyme A synthetase (AMP forming). *Cell Mol Life Sci* **61**, 2020–2030 (2004).
21. Rosenzweig, R. F., Sharp, R. R., Treves, D. S. & Adams, J. Microbial evolution in a simple unstructured environment: genetic differentiation in *Escherichia coli*. *Genetics* **137**, 903–917 (1994).

22. Treves, D. S., Manning, S. & Adams, J. Repeated evolution of an acetate-crossfeeding polymorphism in long-term populations of *Escherichia coli*. *Mol Biol Evol* **15**, 789–797 (1998).
23. Lin, H., Castro, N. M., Bennett, G. N. & San, K.-Y. Acetyl-CoA synthetase overexpression in *Escherichia coli* demonstrates more efficient acetate assimilation and lower acetate accumulation: a potential tool in metabolic engineering. *Appl Microbiol Biotechnol* **71**, 870–874 (2006).
24. Kinnersley, M. A., Holben, W. E. & Rosenzweig, F. E Unibus Plurum: genomic analysis of an experimentally evolved polymorphism in *Escherichia coli*. *PLoS Genet* **5**, e1000713 (2009).
25. Penzlin, A. *et al.* Pipasic: similarity and expression correction for strain-level identification and quantification in metaproteomics. *Bioinformatics* **30**, i149–156 (2014).
26. Ram, R. J. *et al.* Community proteomics of a natural microbial biofilm. *Science* **308**, 1915–1920 (2005).
27. VerBerkmoes, N. C., Denef, V. J., Hettich, R. L. & Banfield, J. F. Systems biology: Functional analysis of natural microbial consortia using community proteomics. *Nat Rev Microbiol* **7**, 196–205 (2009).
28. Zhang, X. *et al.* MetaPro-IQ: a universal metaproteomic approach to studying human and mouse gut microbiota. *Microbiome* **4**, 31 (2016).
29. Abrams, P. A. Character displacement and niche shift analyzed using consumer-resource models of competition. *Theor Popul Biol* **29**, 107–160 (1986).
30. Ispolatov, I. & Doebeli, M. A note on the complexity of evolutionary dynamics in a classic consumer-resource model. *Theor Ecol* **13**, 79–84 (2020).
31. Arthur, R. M. Species Packing, and What Competition Minimizes. *PNAS* **64**, 1369–1371 (1969).
32. You, C. *et al.* Coordination of bacterial proteome with metabolism by cyclic AMP signalling. *Nature* **500**, 301–306 (2013).
33. Basan, M. *et al.* Overflow metabolism in *Escherichia coli* results from efficient proteome allocation. *Nature* **528**, 99–104 (2015).
34. Hardin, G. The competitive exclusion principle. *Science* **131**, 1292–1297 (1960).
35. Tatusov, R. L., Galperin, M. Y., Natale, D. A. & Koonin, E. V. The COG database: a tool for genome-scale analysis of protein functions and evolution. *Nucleic Acids Res* **28**, 33–36 (2000).
36. Tatusov, R. L. *et al.* The COG database: an updated version includes eukaryotes. *BMC Bioinformatics* **4**, 41 (2003).
37. Segata, N. *et al.* Metagenomic microbial community profiling using unique clade-specific marker genes. *Nat Methods* **9**, 811–814 (2012).
38. Li, J. *et al.* An integrated catalog of reference genes in the human gut microbiome. *Nat Biotechnol* **32**, 834–841 (2014).
39. Almeida-Neto, M., Guimarães, P., Guimarães Jr, P. R., Loyola, R. D. & Ulrich, W. A consistent metric for nestedness analysis in ecological systems: reconciling concept and measurement. *Oikos* **117**, 1227–1239 (2008).
40. Sørensen, M. A., Fricke, J. & Pedersen, S. Ribosomal protein S1 is required for translation of most, if not all, natural mRNAs in *Escherichia coli* in vivo<sup>1</sup> Edited by D. Draper. *Journal of Molecular Biology* **280**, 561–569 (1998).

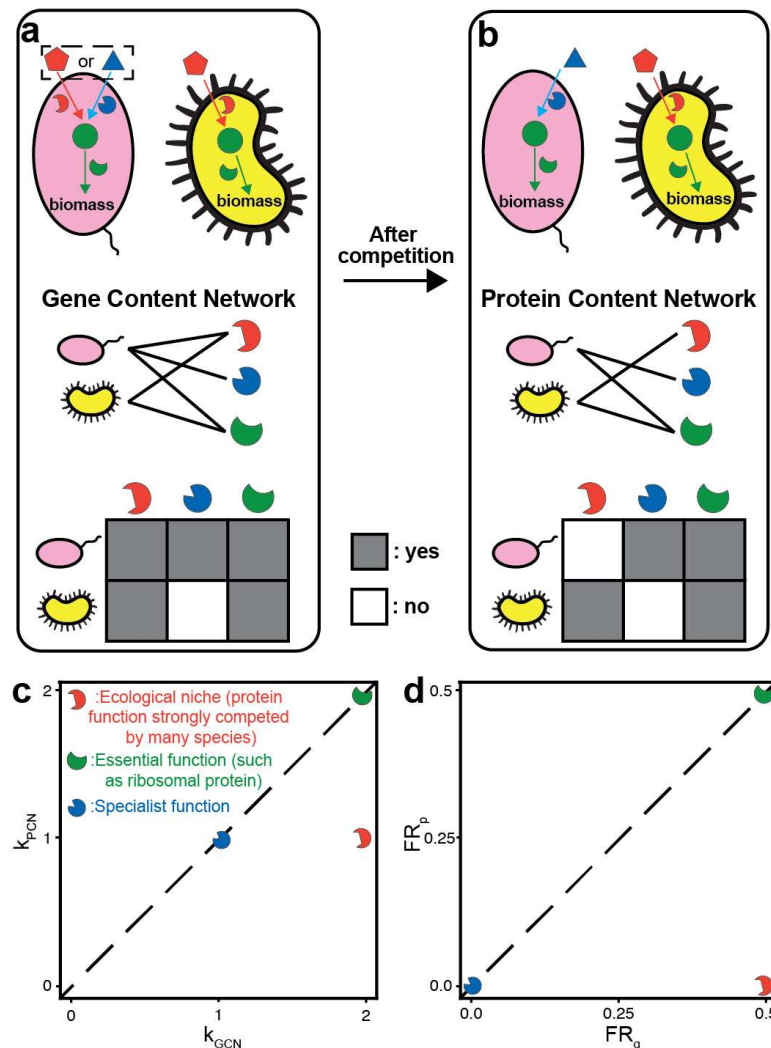
41. Galperin, M. Y., Wolf, Y. I., Garushyants, S. K., Vera Alvarez, R. & Koonin, E. V. Nonessential Ribosomal Proteins in Bacteria and Archaea Identified Using Clusters of Orthologous Genes. *Journal of Bacteriology* **203**, e00058-21 (2021).
42. Shoji, S., Dambacher, C. M., Shajani, Z., Williamson, J. R. & Schultz, P. G. Systematic chromosomal deletion of bacterial ribosomal protein genes. *J Mol Biol* **413**, 751–761 (2011).
43. Dabbs, E. R. Mutants lacking individual ribosomal proteins as a tool to investigate ribosomal properties. *Biochimie* **73**, 639–645 (1991).
44. Baba, T. *et al.* Construction of Escherichia coli K-12 in-frame, single-gene knockout mutants: the Keio collection. *Mol Syst Biol* **2**, 2006.0008 (2006).
45. Steinsiek, S. & Bettenbrock, K. Glucose transport in Escherichia coli mutant strains with defects in sugar transport systems. *J Bacteriol* **194**, 5897–5908 (2012).
46. Fath, M. J. & Kolter, R. ABC transporters: bacterial exporters. *Microbiol Rev* **57**, 995–1017 (1993).
47. Nikaido, H. Maltose transport system of Escherichia coli: an ABC-type transporter. *FEBS Lett* **346**, 55–58 (1994).
48. Yilmaz, B. & Li, H. Gut Microbiota and Iron: The Crucial Actors in Health and Disease. *Pharmaceuticals* **11**, 98 (2018).
49. Seyoum, Y., Baye, K. & Humblot, C. Iron homeostasis in host and gut bacteria – a complex interrelationship. *Gut Microbes* **13**, 1874855 (2021).
50. Cech, T. R. The Ribosome Is a Ribozyme. *Science* **289**, 878–879 (2000).
51. Xue, S. & Barna, M. Specialized ribosomes: a new frontier in gene regulation and organismal biology. *Nat Rev Mol Cell Biol* **13**, 355–369 (2012).
52. Bubunenko, M., Baker, T. & Court, D. L. Essentiality of Ribosomal and Transcription Antitermination Proteins Analyzed by Systematic Gene Replacement in Escherichia coli. *Journal of Bacteriology* **189**, 2844–2853 (2007).
53. Kanehisa, M. A database for post-genome analysis. *Trends Genet* **13**, 375–376 (1997).
54. Kanehisa, M. & Goto, S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res* **28**, 27–30 (2000).
55. Mao, X., Cai, T., Olyarchuk, J. G. & Wei, L. Automated genome annotation and pathway identification using the KEGG Orthology (KO) as a controlled vocabulary. *Bioinformatics* **21**, 3787–3793 (2005).
56. Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a reference resource for gene and protein annotation. *Nucleic Acids Res* **44**, D457-462 (2016).
57. Kotrba, P., Inui, M. & Yukawa, H. Bacterial phosphotransferase system (PTS) in carbohydrate uptake and control of carbon metabolism. *Journal of Bioscience and Bioengineering* **92**, 502–517 (2001).
58. Romano, A. H. & Conway, T. Evolution of carbohydrate metabolic pathways. *Research in Microbiology* **147**, 448–455 (1996).
59. Yan, Y. *Engineering Microbial Metabolism For Chemical Synthesis: Reviews And Perspectives*. (World Scientific, 2017).
60. E, H. G. The multivariate niche. *Cold Spring Harbor Symposia on Quantitative Biology* **22**, 415–421 (1957).
61. Leibold, M. A. The Niche Concept Revisited: Mechanistic Models and Community Context. *Ecology* **76**, 1371–1382 (1995).
62. Holt, R. D. Bringing the Hutchinsonian niche into the 21st century: Ecological and evolutionary perspectives. *PNAS* **106**, 19659–19665 (2009).



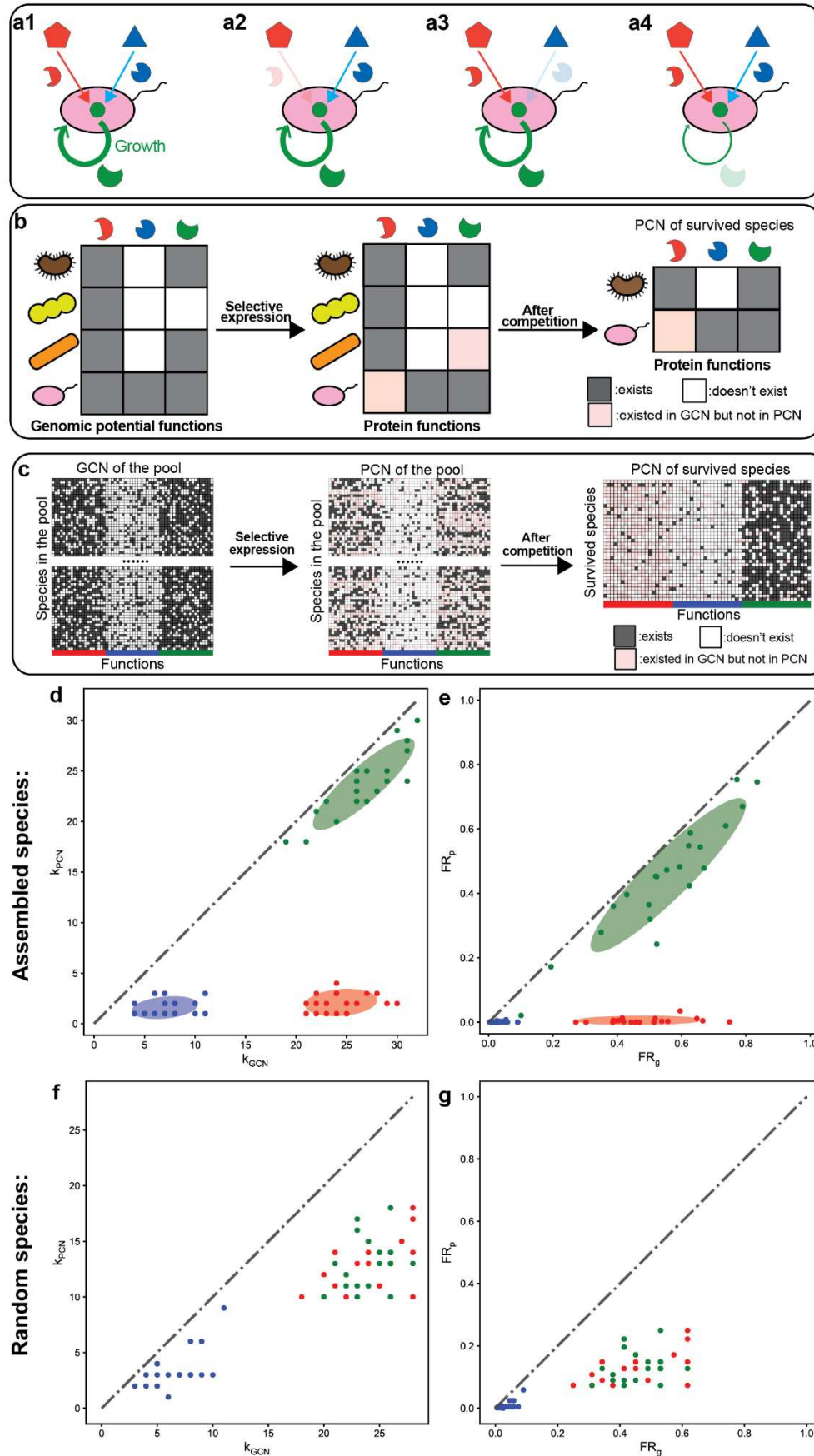
63. Li, L. *et al.* RapidAIM: a culture- and metaproteomics-based Rapid Assay of Individual Microbiome responses to drugs. *Microbiome* **8**, 33 (2020).
64. Li, L. *et al.* An in vitro model maintaining taxon-specific functional activities of the gut microbiome. *Nat Commun* **10**, 4146 (2019).
65. Creskey, M. *et al.* An economic and robust TMT labeling approach for high throughput proteomic and metaproteomic analysis. 2022.07.30.502163 Preprint at <https://doi.org/10.1101/2022.07.30.502163> (2022).
66. Tilman, D. Resource Competition between Plankton Algae: An Experimental and Theoretical Approach. *Ecology* **58**, 338–348 (1977).
67. Wandersman, C. & Delepelaire, P. Bacterial iron sources: from siderophores to hemophores. *Annu Rev Microbiol* **58**, 611–647 (2004).
68. Hibbing, M. E., Fuqua, C., Parsek, M. R. & Peterson, S. B. Bacterial competition: surviving and thriving in the microbial jungle. *Nature Reviews Microbiology* **8**, 15–25 (2010).
69. Smith, R. L. & Smith, T. M. *Elements of ecology*. (Benjamin Cummings, 2003).
70. Hutchinson, G. E. The Paradox of the Plankton. *The American Naturalist* **95**, 137–145 (1961).
71. Marsland, R., Cui, W. & Mehta, P. The Minimum Environmental Perturbation Principle: A New Perspective on Niche Theory. *The American Naturalist* **196**, 291–305 (2020).
72. Rozen, D. E. & Lenski, R. E. Long-Term Experimental Evolution in Escherichia coli. VIII. Dynamics of a Balanced Polymorphism. *Am. Nat.* **155**, 24–35 (2000).
73. Good, B. H., McDonald, M. J., Barrick, J. E., Lenski, R. E. & Desai, M. M. The dynamics of molecular evolution over 60,000 generations. *Nature* **551**, 45–50 (2017).
74. MacArthur, R. Species packing and competitive equilibrium for many species. *Theoretical Population Biology* **1**, 1–11 (1970).
75. Goyal, A. & Maslov, S. Diversity, Stability, and Reproducibility in Stochastically Assembled Microbial Ecosystems. *Phys. Rev. Lett.* **120**, 158102 (2018).
76. Goldford, J. E. *et al.* Emergent simplicity in microbial community assembly. *Science* **361**, 469–474 (2018).
77. Marsland, R. *et al.* Available energy fluxes drive a transition in the diversity, stability, and functional structure of microbial communities. *PLoS Comput Biol* **15**, e1006793 (2019).
78. Li, Z. *et al.* Modeling microbial metabolic trade-offs in a chemostat. *PLoS Comput Biol* **16**, e1008156 (2020).
79. Dubinkina, V., Fridman, Y., Pandey, P. P. & Maslov, S. Multistability and regime shifts in microbial communities explained by competition for essential nutrients. *Elife* **8**, e49720 (2019).
80. Cappelletti, V. *et al.* Dynamic 3D proteomes reveal protein functional alterations at high resolution in situ. *Cell* **184**, 545–559.e22 (2021).
81. Schreiber, F. *et al.* Phenotypic heterogeneity driven by nutrient limitation promotes growth in fluctuating environments. *Nat Microbiol* **1**, 1–7 (2016).
82. Gutierrez-Ríos, R. M. *et al.* Identification of regulatory network topological units coordinating the genome-wide transcriptional response to glucose in Escherichia coli. *BMC Microbiology* **7**, 53 (2007).
83. Goyal, A., Bittleston, L. S., Leventhal, G. E., Lu, L. & Cordero, O. X. Interactions between strains govern the eco-evolutionary dynamics of microbial communities. *Elife* **11**, e74987 (2022).



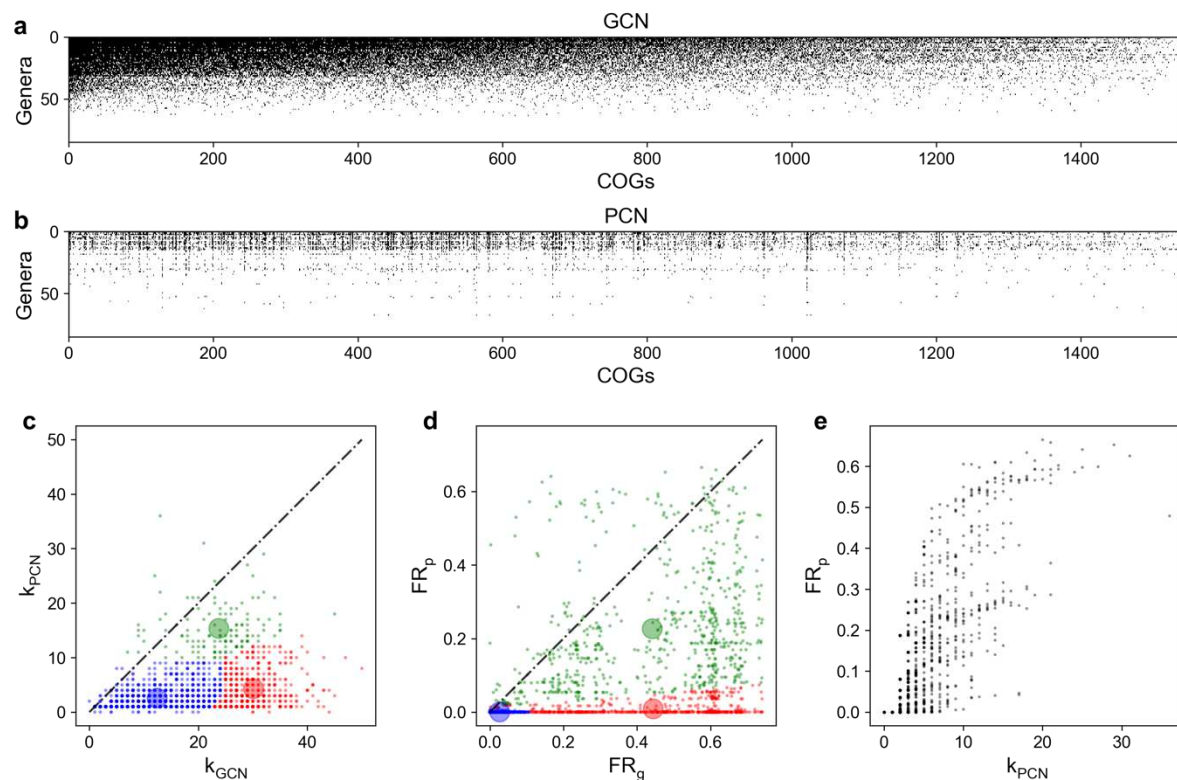
84. Zhang, X. *et al.* Evaluating live microbiota biobanking using an ex vivo microbiome assay and metaproteomics. *Gut Microbes* **14**, 2035658 (2022).
85. Zhang, X. *et al.* Assessing the impact of protein extraction methods for human gut metaproteomics. *Journal of Proteomics* **180**, 120–127 (2018).
86. Perez-Riverol, Y. *et al.* The PRIDE database resources in 2022: a hub for mass spectrometry-based proteomics evidences. *Nucleic Acids Research* **50**, D543–D552 (2022).
87. Tyanova, S., Temu, T. & Cox, J. The MaxQuant computational platform for mass spectrometry-based shotgun proteomics. *Nat Protoc* **11**, 2301–2319 (2016).
88. Cheng, K. *et al.* MetaLab 2.0 Enables Accurate Post-Translational Modifications Profiling in Metaproteomics. *J. Am. Soc. Mass Spectrom.* **31**, 1473–1482 (2020).
89. Huang, T. *et al.* MSstatsTMT: Statistical Detection of Differentially Abundant Proteins in Experiments with Isobaric Labeling and Multiple Mixtures. *Mol Cell Proteomics* **19**, 1706–1723 (2020).
90. Ma, W. *et al.* DreamAI: algorithm for the imputation of proteomics data. 2020.07.21.214205 Preprint at <https://doi.org/10.1101/2020.07.21.214205> (2021).
91. Kanehisa, M., Sato, Y. & Morishima, K. BlastKOALA and GhostKOALA: KEGG Tools for Functional Characterization of Genome and Metagenome Sequences. *J Mol Biol* **428**, 726–731 (2016).
92. Atmar, W. & Patterson, B. D. The measure of order and disorder in the distribution of species in fragmented habitat. *Oecologia* **96**, 373–382 (1993).



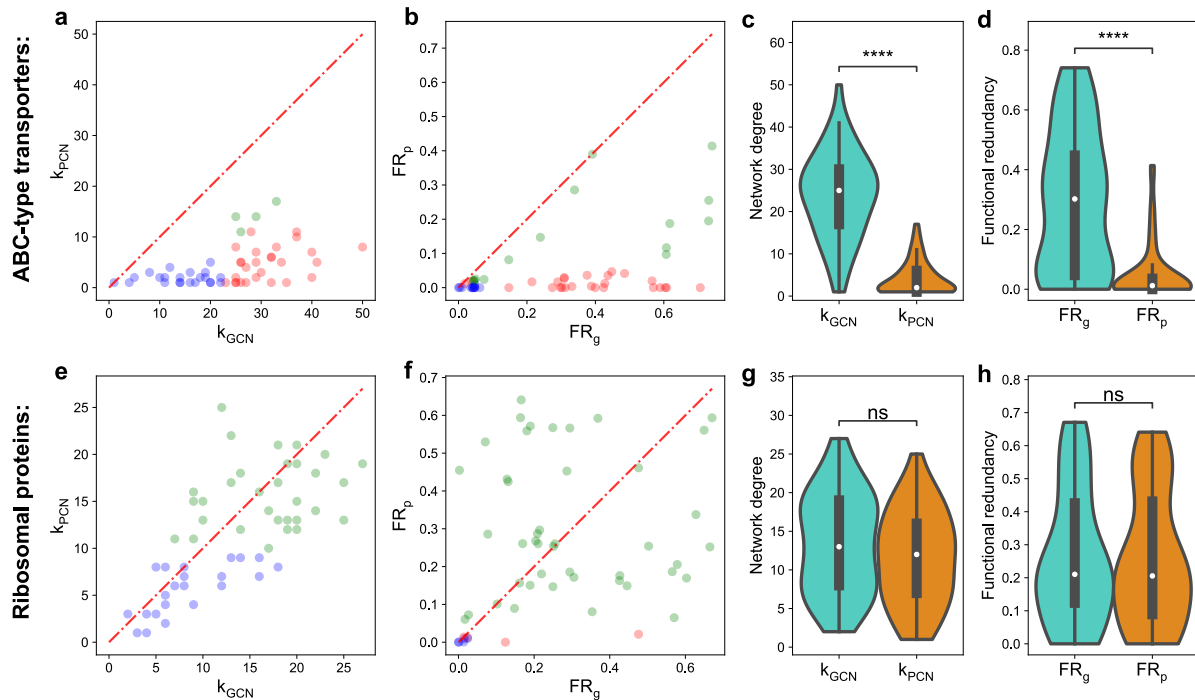
**Figure 1: Protein functions involved in determining ecological niches are postulated to have larger discrepancies between the gene-level functional redundancy  $FR_g$  and protein-level functional redundancy  $FR_p$ .** Here we use a hypothetical example with three representative proteins (3 broken circles with complementary shapes to their substrates) to demonstrate this point. **a**, Schematic of genomic capacity of two microbial taxa (pink oval vs yellow indented oval). Two resources (red pentagon and blue triangle) are externally supplied to the community. The green metabolite can be transformed from the red or blue resource and further utilized in biomass synthesis. The pink taxon has the capacity of converting either supplied resource into the green metabolite (red and blue arrows), while the yellow taxon can only convert the red resource (red arrow). **b**, Schematic of expressed proteins for two microbial taxa after their competition in the same community. After the competition, the reduced resource conflict (represented by the pink taxon choosing the blue resource as the sole one to consume) can promote their coexistence. Gene content network (GCN) and protein content network (PCN) can be used to capture genomic capacity and expressed protein functions for all taxa. Alternatively, this network can be represented as incidence matrices on the bottom (grey areas imply the existence of edges connecting taxa to proteins). **c-d**, The comparison between  $k_{GCN}$  and  $k_{PCN}$  or between  $FR_g$  and  $FR_p$  helps to classify proteins into three protein functional clusters: specialist function, essential function, and niche function. In the calculation of  $FR_g$  and  $FR_p$ , we assume equal abundances of the two species, i.e.,  $p_1 = p_2 = 0.5$ .



**Figure 2: Three protein functional clusters (specialist function, essential function, and niche function) considered in the community assembly model form three distinct clusters when the network degree and functional redundancy are compared between the GCN and PCN in model-generated synthetic data.** **a1-a4**, Three types of functions modeled have different ecological and metabolic roles. The niche function (red proteins) and specialist function (blue proteins) are modeled as abilities to consume externally supplied resources. The role of essential functions (green proteins) is considered as a reduction in the overall growth rate for each missing essential function. **b**, A schematic diagram of the community assembly. Species (ovals and indented ovals) with expressed protein functions selected via the sub-sampling of their genomic capacity. Then all species are co-cultured together to simulate their ecological competition. **c**, A simulation example of the community assembly, and the construction of GCN and PCN for the survived species. **d-e**, The comparison of network degree and functional redundancy respectively based on the GCN and PCN of survived species in the simulation example in panel-c. A Gaussian mixture model with 3 clusters is used to identify 3 protein functional clusters. Ellipses around clusters cover areas one standard deviation away from their means. **f-g**, The comparison of network degree and functional redundancy respectively based on the GCN and PCN of 35 species randomly selected from the 10,000 species in the initial pool. All points/functions are colored red (niche functions), green (essential functions), and blue (specialist functions) according to their types of functions in the model.  $k_{GCN}$  (or  $k_{PCN}$ ) is the network degree of each function in the GCN (or PCN).  $FR_g$  (or  $FR_p$ ) is the functional redundancy of each function on the gene level (or protein level), respectively.

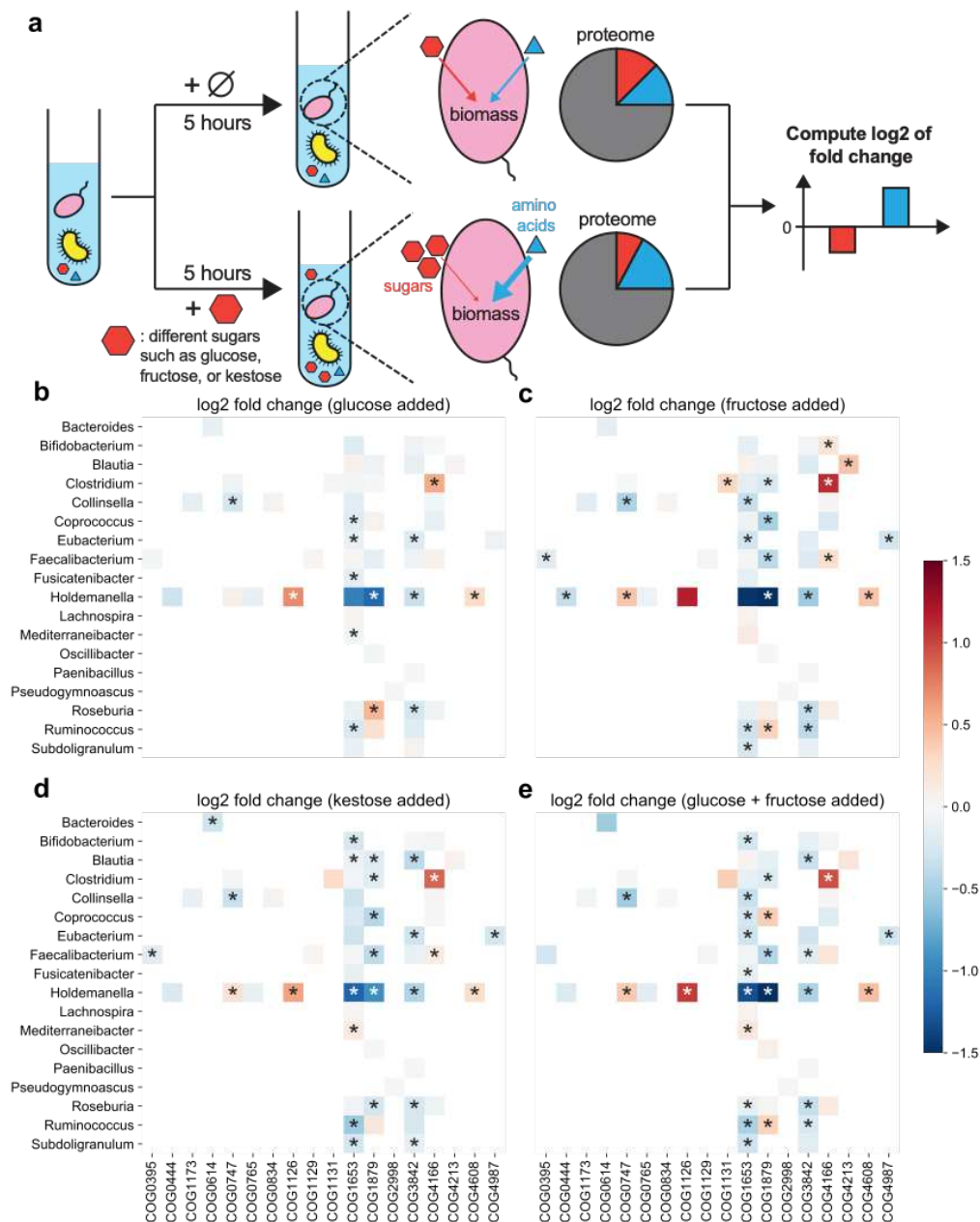


**Figure 3: Real data of the human gut microbiome showing three clusters on the plot that compares  $FR_g$  with  $FR_p$ .** Metagenome and metaproteome of subject HM454 mucosal-luminal interface samples<sup>28</sup> were used to construct GCN and PCN, respectively. **a**, The GCN shows if a genus owns (or doesn't own) a COG as its genomic capacity, which is colored in black (or white). The GCN matrix is ordered to have decreasing network degrees for both genera and COGs. **b**, The PCN shows if a genus expresses (or doesn't express) a COG as its protein function, which is colored in black (or white). The PCN matrix follows the same order as the GCN. **c**, Differences in network degree for most COGs are large.  $k_{GC}$  is the network degree of each COG in the GCN (i.e. the number of genera owning each COG in the GCN).  $k_{PC}$  is the network degree of each COG in the PCN (i.e. the number of genera owning each COG in the PCN). **d**,  $FR_g$  is larger than  $FR_p$  for most COGs. Three clusters with three distinct colors (blue, red, and green) are predicted by the Gaussian mixture model with 3 clusters fitted on synthetic data. The transparent large circles represent centroids of three clusters. **e**, The relationship between  $FR_p$  and network degree of PCN for COGs is not monotonic.



**Figure 4: Comparison of network degree and functional redundancy between the gene and protein level for ABC-type transporters and ribosomal proteins.** **a**, Network degrees in GCN are larger than network degrees in PCN for most ABC-type transporter COGs.  $k_{GCN}$  (or  $k_{PCN}$ ) is the network degree of each COG in the GCN (or PCN). **b**,  $FR_g$  is larger than  $FR_p$  for most ABC-type transporter COGs. **c-d**, The distribution of network degrees and functional redundancies (violin plots and boxplots) for ABC-type transporter COGs show a significantly huge reduction from  $k_{GCN}$  to  $k_{PCN}$  or from  $FR_g$  to  $FR_p$ . **e**, Network degrees in GCN are comparable with that in PCN for most ribosomal protein COGs. **f**,  $FR_g$  is comparable with  $FR_p$  for most ribosomal protein COGs. Points in scatter plots are colored by the same colors used in Fig. 3d. **g-h**, The distribution of network degrees and functional redundancies (violin plots and boxplots) for ribosomal protein COGs show no significant reduction from  $k_{GCN}$  to  $k_{PCN}$  or from  $FR_g$  to  $FR_p$ . In all boxplots, the middle white dot is the median, the lower and upper hinges correspond to the first and third quartiles, and the black line ranges from the  $1.5 \times IQR$  (where  $IQR$  is the interquartile range) below the lower hinge to  $1.5 \times IQR$  above the upper hinge. All violin plots are smoothed by a kernel density estimator and 0 is set as the lower bound. All statistical analyses were performed using the two-sided Mann-Whitney-Wilcoxon U Test with Bonferroni correction between genomic capacity (GCN) and protein functions (PCN). P values obtained from the test is divided into 5 groups: (1)  $p > 0.05$  (ns), (2)  $0.01 < p \leq 0.05$  (\*), (3)  $10^{-3} < p \leq 0.01$  (\*\*), (4)  $10^{-4} < p \leq 10^{-3}$  (\*\*\*), and (5)  $p \leq 10^{-4}$  (\*\*\*\*). Network degree comparison of ABC transporters:  $p = 7.11 \times 10^{-16}$ . Network degree comparison of ribosomal proteins:  $p = 0.10$ . Redundancy comparison of ABC transporters:  $p = 2.19 \times 10^{-11}$ . Redundancy comparison of ribosomal proteins:  $p = 1.00$ .





**Figure 5: Microbes modify their expression for ABC-type transporters to adapt to added sugars.** All heatmaps share the same color bar on the right. **a**, Schematic of in-vitro cultures of a collected human gut microbiome. In the treatment group, one sugar is added to the community. Metaproteomic measurements 5 hours later enable us to compare the intensity of each taxon-specific protein using the log2 fold change of each protein's fraction (i.e. normalized intensity over each genus) from the treatment group divided by that from the control group. Log2 fold changes of ABC-type transporters 5 hours after (b) glucose, (c) fructose, (d) kestose, or (e) glucose and fructose is added.