

# Evidence of cellulose metabolism by the giant panda gut microbiome

Lifeng Zhu<sup>a,1</sup>, Qi Wu<sup>a,1</sup>, Jiayin Dai<sup>a</sup>, Shanning Zhang<sup>b</sup>, and Fuwen Wei<sup>a,2</sup>

<sup>a</sup>Key Laboratory of Animal Ecology and Conservation Biology, Institute of Zoology, Chinese Academy of Sciences, Beijing 100101, China and <sup>b</sup>China Wildlife Conservation Association, Beijing 100714, China

Edited by Rita R. Colwell, University of Maryland, College Park, MD, and approved September 7, 2011 (received for review December 2, 2010)

The giant panda genome codes for all necessary enzymes associated with a carnivorous digestive system but lacks genes for enzymes needed to digest cellulose, the principal component of their bamboo diet. It has been posited that this iconic species must therefore possess microbial symbionts capable of metabolizing cellulose, but these symbionts have remained undetected. Here we examined 5,522 prokaryotic ribosomal RNA gene sequences in wild and captive giant panda fecal samples. We found lower species richness of the panda microbiome than of mammalian microbiomes for herbivores and nonherbivorous carnivores. We detected 13 operational taxonomic units closely related to *Clostridium* groups I and XIVa, both of which contain taxa known to digest cellulose. Seven of these 13 operational taxonomic units were unique to pandas compared with other mammals. Metagenomic analysis using ~37-Mbp contig sequences from gut microbes recovered putative genes coding two cellulose-digesting enzymes and one hemicellulose-digesting enzyme, cellulase,  $\beta$ -glucosidase, and xylan 1,4- $\beta$ -xylosidase, in *Clostridium* group I. Comparing glycoside hydrolase profiles of pandas with those of herbivores and omnivores, we found a moderate abundance of oligosaccharide-degrading enzymes for pandas (36%), close to that for humans (37%), and the lowest abundance of cellulases and endohemicellulases (2%), which may reflect low digestibility of cellulose and hemicellulose in the panda's unique bamboo diet. The presence of putative cellulose-digesting microbes, in combination with adaptations related to feeding, physiology, and morphology, show that giant pandas have evolved a number of traits to overcome the anatomical and physiological challenge of digesting a diet high in fibrous matter.

Access to dietary resources shapes animal evolution (1). Early on, animals lost the ability to synthesize many key compounds, and instead this function is performed by symbionts (2). For example, microbial symbionts assist with extracting nutrients from food and key compounds from the environment, and also synthesize necessary metabolic compounds (1). Gut microbiota share specialized relationships with their hosts, and advances in genomics are revealing the dynamics of these relationships (3). Recent developments in culture-independent methodologies based on large-scale comparative analyses of microbial small-subunit ribosomal RNA genes (16S ribosomal RNA) and metagenomics have revealed the extent of microbial diversity and metabolic potential in greater detail (2–7). These techniques can now be applied to animals that have acquired a profoundly new diet, presenting an opportunity to investigate host physiological and microbial systems in an evolutionary context.

The giant panda (*Ailuropoda melanoleuca*) is well known for dietary oddities: a bamboo specialist within the mammalian order Carnivora possessing a gastrointestinal tract typical of carnivores. It consumes ~12.5 kg of this highly fibrous plant each day (8), but because it lacks the long intestinal tract characteristic of other herbivores, extensive fermentation is not possible (9). Giant pandas digest only ~17% of dry matter consumed (8), and have low digestion coefficients for bamboo hemicelluloses (27%) and celluloses (8%) (9). Indeed, the giant panda genome codes for all necessary enzymes associated with a carnivorous digestive system, but lacks the enzyme homologs needed for

cellulose digestion (10). Although the giant panda can use non-cellulosic material from the bamboo diet using enzymes coded in its own genome, digestion of cellulose and hemicellulose is impossible based on the panda's genetic composition, and must be dependent on gut microbiome. However, previous research using culture methods and small-scale sequencing identified three predominant bacteria from the panda gut—*Escherichia coli*, *Streptococcus*, and Enterobacteriaceae—none of which aids in cellulose digestion (11–13). Thus, an incomplete understanding of the gut microbial ecosystem in this interesting and high-profile species remains because of restrictions in methodology and past reliance on studies of captive animals.

## Results and Discussion

We undertook a large-scale analysis of 16S rRNA gene sequences to profile microbial flora inhabiting the digestive system of giant pandas and used a metagenomic approach based on next-generation de novo sequencing to identify functional attributes encoded in the gut microbiome. A total of 5,636 near-full-length 16S rRNA gene segments were amplified from fecal samples of seven wild and eight captive giant pandas. After exclusion of 74 putative chimeric and 30 chloroplast sequences, 5,522 sequences were retained for analysis. Using a minimum identity of 97% as the threshold for any sequence pair, we identified 85 bacterial operational taxonomic units (OTUs), 14 of which were previously undescribed (Fig. 1A and *SI Appendix*, Table S1). Coverage of 99% was obtained across existing bacterial clone libraries, and we are confident that our dataset presents the most comprehensive assessment of gut microbes in this species based on the rarefaction method in DOTUR (14) (*SI Appendix*, Fig. S1A). The majority of microbes were members of the Firmicutes (62 OTUs, 4,633 sequences, 83.8% of the total of 5,522 sequences) and Proteobacteria (12 OTUs, 871 sequences, 15.8% of the total sequences), with the remainder belonging to the phyla Actinobacteria, Bacteroidetes, Cyanobacteria, and Acidobacteria (Fig. 1 and *SI Appendix*, Fig. S2 and Table S1). Within the Firmicutes, 33 OTUs (60.8% of the total of 5,522 sequences) were members of the class Clostridia and 29 OTUs (23.0% of total sequences) belonged to the class Bacilli. The high proportion of Firmicutes OTUs found in the gut of giant pandas differs from the findings in previous studies attempting to characterize giant panda gut microbes (11–13) and is notably similar

Author contributions: F.W. designed research; L.Z. and S.Z. performed research; L.Z., Q.W., J.D., S.Z., and F.W. analyzed data; and L.Z., Q.W. and F.W. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Freely available online through the PNAS open access option.

Data deposition: The sequences reported in this paper have been deposited in the GenBank database. For a list of accession numbers, see *SI Appendix*. The metagenomic data project ID is 5936 in Integrated Microbial Genomes with Microbiome Samples (IMG/M).

<sup>1</sup>L.Z. and Q.W. contributed equally to this work.

<sup>2</sup>To whom correspondence should be addressed. E-mail: weifw@ioz.ac.cn.

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.1017956108/-DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1017956108/-DCSupplemental).





**Table 1. Inventory of putative GHs identified in the giant panda gut microbiome**

CAZy family	Pfam HMM name	Pfam accession	Known activities	Giant panda gut community
GH1	Glyco_hydro_1	PF00232	$\beta$ -Glucosidase, $\beta$ -galactosidase, $\beta$ -mannosidase, others	101
GH2	Glyco_hydro_2_C	PF00703	$\beta$ -Galactosidase, $\beta$ -mannosidase, others	1
GH3	Glyco_hydro_3	PF00933	$\beta$ -1,4-Glucosidase, $\beta$ -1,4-xylosidase, $\beta$ -1,3-glucosidase, $\alpha$ -l-arabinofuranosidase, others	18
GH4	Glyco_hydro_4	PF02056	$\alpha$ -Glucosidase, $\alpha$ -galactosidase, $\alpha$ -glucuronidase, others	6
GH5	Cellulase	PF00150	Cellulase, $\beta$ -1,4-endoglucanase, $\beta$ -1,3-glucosidase, $\beta$ -1,4-endoxylanase, $\beta$ -1,4-endomannanase, others	3
GH8	Glyco_hydro_8	PF01270	Cellulase, $\beta$ -1,3-glucosidase, $\beta$ -1,4-endoxylanase, $\beta$ -1,4-endomannanase, others	2
GH10	Glyco_hydro_10	PF00331	Xylanase, $\beta$ -1,3-endoxylanase	2
GH13	Alpha-amylase	PF00128	$\alpha$ -Amylase, catalytic domain, and related enzymes	35
GH16	Glyco_hydro_16	PF00722	$\beta$ -1,3(4)-Endoglucanase, others	9
GH18	Glyco_hydro_18	PF00704	Chitinase, endo- $\beta$ - <i>N</i> -acetylglucosaminidase, noncatalytic proteins	1
GH20	Glyco_hydro_20	PF00728	$\beta$ -Hexosaminidase, lacto- <i>N</i> -biosidase	7
GH23	SLT	PF01464	G-type lysozyme, peptidoglycan lytic transglycosylase	20
GH24	Phage_lysozyme	PF00959	Lysozyme	1
GH25	Glyco_hydro_25	PF01183	Lysozyme	9
GH27	Melibiose	PF02065	$\alpha$ -Galactosidase, $\alpha$ - <i>N</i> -acetylgalactosaminidase, isomalto-dextranase	36
GH28	Glyco_hydro_28	PF00295	Polygalacturonase, rhamnogalacturonase, others	1
GH29	Alpha_L_fucos	PF01120	$\alpha$ -L-fucosidase	1
GH31	Glyco_hydro_31	PF01055	$\alpha$ -Glucosidase, $\alpha$ -xylosidase, others	9
GH32	Glyco_hydro_32N	PF00251	Levanase, invertase, others	13
GH35	Glyco_hydro_35	PF01301	$\beta$ -Galactosidase	4
GH36	No Pfam	No Pfam	$\alpha$ -Galactosidase, $\alpha$ - <i>N</i> -acetylgalactosaminidase	33
GH37	Trehaalse	PF01204	$\alpha$ , $\alpha$ -trehalase	2
GH38	Glyco_hydro_38	PF01074	$\alpha$ -Mannosidase	10
GH39	Glyco_hydro_39	PF01229	$\beta$ -Xylosidase, $\alpha$ -L-iduronidase	9
GH42	Glyco_hydro_42	PF02449	$\beta$ -Galactosidase	18
GH51	No Pfam	No Pfam	Endoglucanase, $\alpha$ -L-arabinofuranosidase	5
GH55	No Pfam	No Pfam	Exo-1,3-glucanase, endo-1,3-glucanase	1
GH58	No Pfam	No Pfam	Endo- <i>N</i> -acetylneuraminidase or endo-sialidase	1
GH65	Glyco_hydro_65m	PF03632	Trehalase, maltose phosphorylase, trehalose phosphorylase	1
GH67	Glyco_hydro_67M	PF03648	$\alpha$ -Glucuronidase, others	2
GH70	Glyco_hydro_70	PF02324	Dextranucrase, alternansucrase	9
GH72	Glyco_hydro_72	PF03198	$\beta$ -1,3-glucanosyltransglycosylase	3
GH73	Glucosaminidase	PF01832	Peptidoglycan hydrolase with endo- $\beta$ - <i>N</i> -acetylglucosaminidase specificity	28
GH77	Glyco_hydro_77	PF02446	4- $\alpha$ -Glucanotransferase, amyloamaltase	12
GH78	Bac_rhamnosid	PF05592	$\alpha$ -L-Rhamnosidase	2
GH85	Glyco_hydro_85	PF03644	Endo- $\beta$ - <i>N</i> -acetylglucosaminidase	4
GH88	Glyco_hydro_88	PF07470	D-4,5 Unsaturated $\beta$ -glucuronyl hydrolase	1
GH89	NAGLU	PF05089	Alpha- <i>N</i> -acetylglucosaminidase	2
GH90	No Pfam	No Pfam	Endorhamnosidase	1
GH94	No Pfam	No Pfam	Cellulose phosphorylase, chitobiose phosphorylase, cellodextrin phosphorylase	7
GH95	No Pfam	No Pfam	$\alpha$ -L-Fucosidase	4
GH101	No Pfam	No Pfam	Endo- $\alpha$ - <i>N</i> -acetylgalactosaminidase	6
GH104	Phage_lysozyme	PF00959	Peptidoglycan lytic transglycosylase	1
GH109	No Pfam	No Pfam	$\alpha$ - <i>N</i> -acetylgalactosaminidase	7

CAZy family and known activities are from the Carbohydrate Active Enzymes database (<http://www.cazy.org/>) (32). Pfam HMM-based sequences name and Pfam accession are from the Pfam database (<http://pfam.sanger.ac.uk/>) (33), a large collection of protein families, each represented by multiple sequence alignments and hidden Markov models. Giant panda gut community: Number of predicted genes match the Pfam domains or "non-Pfam" representative sequences.

were as follows: 95 °C for 5 min, followed by 35 cycles of 95 °C for 40 s, 55 °C for 45 s, and 72 °C for 2 min, with a final extension period of 10 min at 72 °C. Replicate PCRs were pooled and gel-purified with the Wizard SV Gel and PCR Clean-Up System (Promega). Amplicons were cloned into PMD18-T vector (TaKaRa) and transformed into *E. coli* TOP10 (Tiangen). Approximately 192–900 colonies from each sample PCR product were chosen at random. Plasmid inserts were sequenced bidirectionally using BigDye Terminator (Applied Biosystems) and vector-specific primers (M13F: 5'-GTAAACGACGGCCAG-3';

M13R: 5'-CAGGAAACAGCTATGAC-3'). Sequences were analyzed with DNASTAR (DNASstar) and trimmed to remove vector sequences. After trimming and adjusting for quality values, the average single sequence read length was ~700 nucleotides. Bidirectional sequence reads of clone inserts provided near-full-length 16S rRNA gene sequences of ~1,400 bp.

**Chimera Checking.** Each sequence was edited manually in conjunction with its chromatogram and secondary structure information. Only unambiguous

nucleotide positions were included in the analysis, and primer sequences were excluded. A multiple sequence alignment was generated with the NAST online tool (27), and chimeras were identified with Bellerophon version 3 (28), implemented at the Greengenes Web site (<http://greengenes.lbl.gov>) with the following (default) parameters: Sequences were compared with others within the same host species and with the Greengenes Core Set, identity to the core set was set to 97%, the match length to sequence threshold was set to 1,050 bp and 1,150 bp, respectively, the window size was set to 300, the count of similar sequences to search for each window was 7, the parent-to-fragment ratio was 90%, and the divergence ratio threshold was set at 1.1.

**Determining OTU and Taxonomy Assignments.** *Determining OTU for each sample.* OTUs were defined as terminal nodes in phylogenetic analysis. Sequences remaining after the chimeras were checked for each individual were aligned using CLUSTAL in MEGA4 (29). OTU determination was performed according to DOTUR analysis (14). In DOTUR, sequences were grouped into OTU using distance-similarity matrices (by PHYLIP; <http://evolution.genetics.washington.edu/phylip.html>), such that the least similar pair within the OTU shared at least 97% identity.

*Determining OTUs for all samples.* To increase validity and decrease computing time, we determined all OTUs and shared OTUs in the sequences for all samples in two steps. First, we selected one wild and one captive sample sequence as the central library. We combined each sample sequence with this central library independently, and then determined the shared OTUs between the central libraries with each sample. Second, we combined unshared sequences for each sample into one new library, then used the foregoing method to determine the OTUs in this remaining large library. We chose a representative sequence for each OTU in each of the 97% identity groups at random, then used the BLAST method to find the closest GenBank neighbor for these representative sequences. Each representative sequence for each OTU was submitted to GenBank and granted an accession number. The result of EzTaxon (30), which was also used to determine the closest taxon for each OTU, served as the reference. Further, each representative sequence for each OTU with <95% identity to any GenBank sequence was defined as a previously undescribed OTU.

**Phylogenetic Analysis.** Based on the foregoing taxon assignment results, we downloaded the representative of the 16S rRNA gene sequences of the bacteria from the National Center for Biotechnology Information (NCBI) Nucleotide Database, combining all OTUs to construct the neighbor-joining phylogenetic tree (1,000 bootstraps) using MEGA4 (29). We selected the representative sequence for each OTU at random. We used this tree to determine the phylogenetic position for each out, and used the representative sequence on each OTU to construct the neighbor-joining phylogenetic tree (1,000 bootstraps) with MEGA4 (29).

**Relationship Between All OTUs and Known *Clostridium* Clusters.** We downloaded sequences from GenBank according to known *Clostridium* clusters (15). We then combined these sequences with all of the OTUs in this study and constructed the neighbor-joining phylogenetic tree (1,000 bootstraps) using MEGA4 (29). Based on phylogenetic information (in the same clade), we determined the relationships between all OTUs and known *Clostridium* clusters.

**Estimation of Microbial Diversity.** We calculated the Good coverage estimate as  $[1 - (n/N)] \times 100$ , where  $n$  is the number of singleton sequences and  $N$  is the total number of sequences for the sample analyzed (3). We used the rarefaction method in DOTUR (14) to explore the diversity of our clone libraries and found that the number of observed OTUs increased with additional sampling effort; however, the 99% coverage obtained over existing bacterial clone libraries indicated that only 1 new OTU would be expected for every 100 additional sequenced clones. Species richness was defined as the number of gut microbe species detected in a sample. DOTUR was used to calculate various diversity indices along with species richness. The classic Shannon–Weaver diversity indices (with 95% confidence interval) of each sample were calculated as described previously (14).

**Metagenomics Analysis: DNA Extraction, DNA Library Construction, and Sequencing.** DNA was extracted from fecal samples of three individual wild giant pandas using the Qiagen QIAamp DNA Stool Mini Kit according to the protocol for isolation of DNA for pathogen detection. Sequencing and general data analyses were performed by BGI (Shenzhen, China). DNA library construction and sequencing followed BGI's previous work on human gut

microbe metagenomic sequencing (7). A library with 350-bp clone insert size was constructed for our sample. Approximately 83 million high-quality reads, with read lengths of 90 bp, were generated for the samples; thus, the total data volume of high-quality reads was nearly 7.5 Gbp (*SI Appendix, Table S3*).

**Public Data Use.** Public data used in the metagenomic analysis included the human gut contig set (version December 10, 2009), the NCBI sequenced bacteria genomes database (version January 13, 2010), the NCBI sequenced fungi genomes database (version March 17, 2010), the NCBI sequenced protozoa genomes database (version July 25, 2007), the Ribosomal Database Project (version July 14, 2010), and the integrated NCBI-NR database (version March 2011).

**Illumina Genome Analyzer (GA) Short Reads de novo Assembly.** We compared the raw short reads with giant panda genome data to remove the host sequence. The clean reads thus obtained were assembled to obtain long contig sequences by the SOAPdenovo assembler (22), as used in human gut microbe metagenomic analyses. We tried different Kmer frequencies to obtain different assembly results, and used N50 lengths to access the best assembly result. The longest contig length and highest read utilization rate were obtained in Kmer55 for W1, Kmer31 for W2, and Kmer53 for W5. Thus, we used the contigs of Kmer55 for W1, Kmer31 for W2, and Kmer53 for W5 as the final assembly result (*SI Appendix, Table S4*).

**Gene Prediction and Taxonomic Assignment.** We used the assembly contig sequences and applied MetaGene software, with only ORFs longer than 100 bp preserved. We translated the ORFs into protein sequences using NCBI Genetic Code 11. We carried out BLASTP (31) alignment to query the predicted protein sequences against the integrated NR protein database. For each predicted gene, hits with E-values  $>1 \times 10^{-5}$  were filtered. Then a significant-matches set was retained to distinguish taxonomic groups, which were defined for hits with E-values  $<10$  times the top hit E-value. Next, the LCA-based algorithm implemented in MEGAN (24) was introduced to determine the taxonomic level of each gene. The LCA-based algorithm assigns genes to taxa so that the taxonomic level of the assigned taxon reflects the level of conservation of the gene. For example, if the hits in the significant matches set belonged to one species, then the predicted gene was considered conserved within the species and assigned to the species. If the hits belonged to several species within one genus, then the predicted gene was assigned to the genus, which was considered the LCA of the predicted gene. If the hits belonged to more species from different genera and all of the genera belonged to one family (or higher taxonomic unit), then the predicted gene was assigned to the family, which was then considered the LCA. The remaining higher taxa can be determined in the same manner.

**Gene Functional Classification.** We performed predicted gene functional classification by querying protein sequences of the genes against the eggNOG database (an integration of the COG and KOG databases) and the KEGG database using BLASTP with E-values  $<1 \times 10^{-5}$ . Genes were annotated as a function of the NOG or KEGG homologs with the lowest E-value. In the COG database, genes were classified into COG categories, whereas in KEGG, genes were assigned to KEGG pathways.

**Glycoside Hydrolase (GH) Family Annotation and Analysis.** We performed database searches for GHs using HMMER HMM-based sequences search with Pfam hidden Markov models (Pfam v25.0 and HMMER v3.0) as described previously (20, 25). We used the Pfam\_ls HMMs to find complete matches with the family by global alignment. All hits with E-values  $<10^{-4}$  were counted, and their sequences were analyzed further. GHs were named in accordance with the CAZy nomenclature scheme (32). For those GH families for which there is currently no Pfam HMM-based sequences, the representative sequences selected from the CAZy Web site were used in BLAST searches of the metagenomic data to identify these GH families using an E-value cutoff of  $10^{-6}$  (20, 25).

**ACKNOWLEDGMENTS.** We thank Drs. Chaodong Zhu and Hengling Cui for assistance with data analysis, Yonggang Nie for assistance in the field, the anonymous reviewers for their constructive comments, Prof. James Elser for the English polishing and giant panda breeding centers for sampling assistance. This study was supported by the National Basic Research Program of China (973 Program; 2007CB411600), the National Natural Science Foundation of China (30830020 and 31070331), and the Knowledge Innovation Program of the Chinese Academy of Sciences (KSCX2-EW-Z-4).

