

# Corrections

## MEDICAL SCIENCES

Correction for “Combined TRPC3 and TRPC6 blockade by selective small-molecule or genetic deletion inhibits pathological cardiac hypertrophy,” by Kinya Seo, Peter P. Rainer, Virginia Shalkey Hahn, Dong-ik Lee, Su-Hyun Jo, Asger Andersen, Ting Liu, Xiaoping Xu, Robert N. Willette, John J. Lepore, Joseph P. Marino, Jr., Lutz Birnbaumer, Christine G. Schnackenberg, and David A. Kass, which appeared in issue 4, January 28, 2014, of *Proc Natl Acad Sci USA* (111:1551–1556; first published January 22, 2014; 10.1073/pnas.1308963111).

The authors note that on page 1555, left column, fourth full paragraph, line 2 “example 17” should instead appear as “example 15.”

[www.pnas.org/cgi/doi/10.1073/pnas.1405263111](http://www.pnas.org/cgi/doi/10.1073/pnas.1405263111)

## MICROBIOLOGY

Correction for “Dengue virus envelope protein domain I/II hinge determines long-lived serotype-specific dengue immunity,” by William B. Messer, Ruklanthi de Alwis, Boyd L. Yount, Scott R. Royal, Jeremy P. Huynh, Scott A. Smith, James E. Crowe, Jr., Benjamin J. Doranz, Kristen M. Kahle, Jennifer M. Pfaff, Laura J. White, Carlos A. Sariol, Aravinda M. de Silva, and Ralph S. Baric, which appeared in issue 5, February 4, 2014, of *Proc Natl Acad Sci USA* (111:1939–1944; first published January 2, 2014; 10.1073/pnas.1317350111).

The authors note that the following statement should be added to the Acknowledgments: “This research was also supported by the Sunlin and Priscilla Chou Foundation (W.B.M).”

[www.pnas.org/cgi/doi/10.1073/pnas.1405351111](http://www.pnas.org/cgi/doi/10.1073/pnas.1405351111)

## ECOLOGY

Correction for “Tackling soil diversity with the assembly of large, complex metagenomes,” by Adina Chuang Howe, Janet K. Jansson, Stephanie A. Malfatti, Susannah G. Tringe, James M. Tiedje, and C. Titus Brown, which appeared in issue 13, April 1, 2014, of *Proc Natl Acad Sci USA* (111:4904–4909; first published March 14, 2014; 10.1073/pnas.1402564111).

The authors note that the accession number 4504979.3 (Iowa corn) should instead appear as [4504797.3](https://doi.org/10.1073/pnas.1405719111) (Iowa corn).

[www.pnas.org/cgi/doi/10.1073/pnas.1405719111](http://www.pnas.org/cgi/doi/10.1073/pnas.1405719111)

## CELL BIOLOGY

Correction for “Sel1L is indispensable for mammalian endoplasmic reticulum-associated degradation, endoplasmic reticulum homeostasis, and survival,” by Shengyi Sun, Guojun Shi, Xuemei Han, Adam B. Francisco, Yewei Ji, Nuno Mendonça, Xiaojing Liu, Jason W. Locasale, Kenneth W. Simpson, Gerald E. Duhamel, Sander Kersten, John R. Yates III, Qiaoming Long, and Ling Qi, which appeared in issue 5, February 4, 2014, of *Proc Natl Acad Sci USA* (111:E582–E591; first published January 22, 2014; 10.1073/pnas.1318114111).

The authors note that on page E590, left column, first paragraph, line 1, “JAX 004781” should instead appear as “JAX 004682.”

[www.pnas.org/cgi/doi/10.1073/pnas.1405563111](http://www.pnas.org/cgi/doi/10.1073/pnas.1405563111)

# Tackling soil diversity with the assembly of large, complex metagenomes

Adina Chuang Howe<sup>a,b,1</sup>, Janet K. Jansson<sup>c,d</sup>, Stephanie A. Malfatti<sup>c</sup>, Susannah G. Tringe<sup>c</sup>, James M. Tiedje<sup>a,b,1</sup>, and C. Titus Brown<sup>a,e</sup>

Departments of <sup>a</sup>Microbiology and Molecular Genetics and <sup>e</sup>Computer Science and Engineering, Michigan State University, East Lansing, MI 48824; <sup>b</sup>Department of Plant, Soil, and Microbial Sciences, Michigan State University, East Lansing, MI 48824; <sup>c</sup>Department of Energy Joint Genome Institute, Walnut Creek, CA 94598; and <sup>d</sup>Earth Sciences Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720

Contributed by James M. Tiedje, February 11, 2014 (sent for review December 10, 2013)

**The large volumes of sequencing data required to sample deeply the microbial communities of complex environments pose new challenges to sequence analysis. De novo metagenomic assembly effectively reduces the total amount of data to be analyzed but requires substantial computational resources. We combine two preassembly filtering approaches—digital normalization and partitioning—to generate previously intractable large metagenome assemblies. Using a human-gut mock community dataset, we demonstrate that these methods result in assemblies nearly identical to assemblies from unprocessed data. We then assemble two large soil metagenomes totaling 398 billion bp (equivalent to 88,000 *Escherichia coli* genomes) from matched Iowa corn and native prairie soils. The resulting assembled contigs could be used to identify molecular interactions and reaction networks of known metabolic pathways using the Kyoto Encyclopedia of Genes and Genomes Orthology database. Nonetheless, more than 60% of predicted proteins in assemblies could not be annotated against known databases. Many of these unknown proteins were abundant in both corn and prairie soils, highlighting the benefits of assembly for the discovery and characterization of novelty in soil biodiversity. Moreover, 80% of the sequencing data could not be assembled because of low coverage, suggesting that considerably more sequencing data are needed to characterize the functional content of soil.**

**C**omplex microbial communities operate at the heart of many crucial terrestrial, aquatic, and host-associated processes, providing critical ecosystem functionality that underpins much of biology (1–7). DNA sequencing has begun to reveal the enormous biological diversity and heterogeneity within these systems, making them difficult to study in situ (2, 4, 5). With ultradeep sequencing, we now have unprecedented access to even the rare species in these environments. However, in complex environments such as soil [where an estimated 50 Tbp is required to sample a gram adequately (8)], converting these large volumes of sequencing data to biologically useful information remains a major challenge.

As the sizes of sequencing datasets grow at an exponential rate, significant computational resources for data storage and analysis are required. A single metagenomic project can readily generate as much or more data than is in global reference databases; for example, a human-gut metagenome sample containing 578 Gbp [ERA000116 (5)], produced more than twice the data in the National Center for Biotechnology Information (NCBI) RefSeq (Release 56) database. In its simplest form, these data (millions to billions of short reads) are error prone and contain only minimal signal for homology searches, limiting direct annotation approaches against reference databases (9). Furthermore, in systems where little of the microbial diversity has been characterized, these annotation approaches are challenged by a lack of reference genomes, and more than half of identified genes share little or no similarity to any experimentally studied genes (1, 5).

Consequently, investigators of environmental metagenomic datasets are confronted by overwhelming volumes of data for

which they have neither the computational resources nor effective bioinformatics tools (because of short read lengths or a lack of reference genomes) to analyze efficiently. De novo assembly of sequence data offers several advantages for analyzing metagenomic datasets. It provides improved accuracy of sequences by removing most random sequencing errors and results in longer and more specific contigs than found in unassembled sequencing reads (10). Furthermore, assembly significantly reduces the total volume of data required for downstream analysis (e.g., gene annotation). Also, de novo assembly does not rely on the existence of reference genomes, thus allowing the discovery of novel genomic elements. The main challenge for metagenomic applications of de novo assembly is that current assembly tools do not scale to the high diversity and large volume of metagenomic data. Metagenomes from rumen, human gut, and permafrost soil sequencing could be assembled only by discarding low-abundance sequences before assembly (2, 4, 5). Although many metagenome-specific assemblers have been developed recently for the assembly of low-complexity communities, they cannot work with the volume of reads necessary to achieve high coverage for extremely diverse environmental metagenomes (10–12).

Here, we apply two preassembly read-filtering strategies, digital normalization and partitioning, that together provide a general strategy for scaling and improving metagenome assembly for large, complex datasets (e.g., billions of reads). Digital

## Significance

**Investigations of complex environments rely on large volumes of sequence data to adequately sample the genetic diversity of a microbial community. The assembly of short-read data into longer, more interpretable sequence currently is not possible for much of the research community because it requires specialized computational facilities. We present approaches that make de novo assembly of complex metagenomes more accessible. These approaches scale data size with community richness and subdivide the data into tractable subsets representing individual species. We applied these methods toward the assembly of two large soil metagenomes to identify important metagenomic references and show that considerably more data are needed to study the terrestrial microbiome comprehensively.**

Author contributions: A.C.H., J.K.J., S.G.T., J.M.T., and C.T.B. designed research; A.C.H. performed research; A.C.H., S.A.M., and C.T.B. contributed new reagents/analytic tools; A.C.H., J.M.T., and C.T.B. analyzed data; and A.C.H., J.M.T., and C.T.B. wrote the paper.

The authors declare no conflict of interest.

Freely available online through the PNAS open access option.

Data deposition: The sequences discussed in this paper have been deposited in the MG-RAST database (<http://metagenomics.anl.gov/>). The unassembled reads can be found as accession nos. [4539514.3–4539528.3](#) (Iowa corn) and [4539571.3–4539594.3](#) (Iowa prairie). The assembled metagenomes can be found as accession nos. [4504979.3](#) (Iowa corn) and [4504798.3](#) (Iowa prairie).

<sup>1</sup>To whom correspondence may be addressed. E-mail: [howead@msu.edu](mailto:howead@msu.edu) or [tiedje@msu.edu](mailto:tiedje@msu.edu).

This article contains supporting information online at [www.pnas.org/lookup/suppl/doi:10.1073/pnas.1402564111/-DCSupplemental](http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.1402564111/-DCSupplemental).

normalization reduces the size of the dataset by setting aside reads from high-coverage regions and results in more uniform sequence coverage overall; digital normalization has been used previously for both genome and mRNA-seq assembly (13). We apply digital normalization to metagenomes to scale assembly by sample richness rather than by diversity. Further, we demonstrate that digital normalization combined with a partitioning approach to separate reads based on transitive connectivity (e.g., grouping reads with sequencing overlap) (14) can be applied to complex metagenomes. Because these approaches have yet to be applied to environmental metagenomes for which the true content is unknown, we first evaluated this strategy on a human-gut mock community (HGMC) dataset containing 21 known genomes and found that these methods result in assemblies that are nearly identical to assemblies from the unprocessed HGMC dataset. Moreover, we show that partitioning separates most reads into species-level bins, providing an alternative to abundance-based and k-mer approaches to species clustering. We next used these approaches to assemble two previously intractable metagenomes from matched soils, 100-y-cultivated Iowa agricultural corn soil and native Iowa prairie soil. We use the resulting assemblies to evaluate our ability both to sample and to characterize small, 3- to 6-g soil samples and their associated functional diversity. Even with 300 Gbp of data, we are unable to achieve deep coverage of the majority of organisms in the sample, highlighting the need for more extensive sequencing.

## Results

**Normalization Results in Similar Assemblies with Minimal Loss of Information.** The HGMC dataset contains sequences from mixed DNA from isolates at varying abundances ranging from fourfold to 2,000-fold sequencing coverage using the Illumina sequencing platform (Table S1). We evaluated our ability to describe the original HGMC genomes and to estimate the abundances of these genomes from our filtered assembly as compared with the unfiltered, original assembly (Fig. 1, Assembly I and Assembly II).

After sequencing, the mock metagenome (unassembled) encompassed a total of 93% of the genomic content of the reference genomes (Fig. S1). After digital normalization, reads were removed based on their coverage within the dataset (Materials and Methods), resulting in a total of 5.9 million reads (40% of the total reads) from the original HGMC dataset (Table 1) with coverage of 91% of the reference genomes (recovery per genome in Fig. S1). The resulting assembly of filtered HGMC reads (normalized) was compared with the assembly of all original reads, evaluating the recovery of reference genomes and the length distribution of assembled contigs for each reference. Using the Velvet assembler (15), we recovered 43% and 44% of the reference genomes in the original and filtered assemblies, respectively. The assembly of the original dataset contained 29,063 contigs and 38 million bp; the filtered assembly contained 30,082 contigs and 35 million bp (Table 3). Comparable recoveries of references between original and filtered datasets also were obtained with other assemblers [SOAPdenovo (16) and

Meta-ITBA (17)]. Overall, the unfiltered and filtered assemblies were very similar, sharing 95% of genomic content (Table S2), and the distributions of contig lengths in unfiltered and filtered assemblies also were comparable. For the large majority of genomes, the filtered assembly recovered similar fractions of each reference. In genomes with lower coverage, such as NC\_003112.2 and NC\_006085.1, improved assemblies from normalization were observed. In genomes with high sequencing coverage, such as the plasmids NC\_005008.1, NC\_005007.1, and NC\_005003.1, the unfiltered assembly recovered significantly more of the original sequence (Table S1).

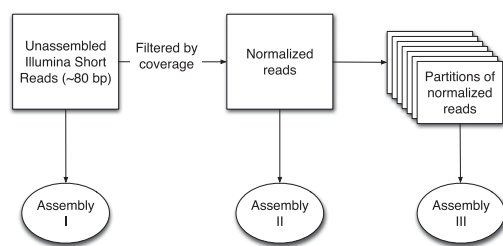
To understand the representation of genes and genomes in the metagenome, we evaluated our ability to estimate genome abundance in the HGMC metagenome and both unfiltered and filtered assemblies. Abundance was estimated through the alignment of unassembled reads to both the known reference genomes or assembled contigs (experimentally hypothesized references). Sequencing coverage was determined as the median base-pair coverage of all aligned reads. For assembled contigs with a coverage greater than 5, the majority of reads that could be aligned to contigs also were mapped to reference genomes (Fig. S2). Below this threshold, reads were mapped to reference genomes but were less likely to be associated with assembled contigs. When the unfiltered and filtered assemblies were compared, the estimated abundance of the HGMC genomes from the filtered assembly were significantly closer to abundances predicted from reference genomes ( $n = 28,652$ ;  $P = 0.032$ ; see SI Materials and Methods).

**Partitioning Separates Most Reads by Species.** To subdivide large metagenomic datasets, we next partitioned the normalized dataset based on De Bruijn graph connectivity. This approach should separate disconnected sequences from distinct species and allow the assembly of each partition independently (Fig. 1, Assembly III). Notably, conserved regions shared by multiple genomes (e.g., 16S rRNA genes) may be connected within a single partition; we examined these partitions through the HGMC dataset. Overall, it was partitioned into 85,818 disconnected partitions containing a total of 9 million reads. Among these, only 2,359 (2.7%) of the partitions contained reads originating from more than one genome, indicating that partitioning separated the large majority of reads from distinct genomes. Partitioning had minimal effects on the assembly of a mock metagenome; the HGMC assemblies of the unpartitioned and partitioned dataset were very similar, sharing 99% identical genomic content.

The number of partitions in the mock metagenome depends largely on the sequencing coverage of its content. In general, reference genomes with high sequence coverage were associated with fewer partitions; a total of 112 partitions contained reads from high-abundance reference genomes (coverage above 25), whereas 2,771 partitions were associated with lower-abundance genomes (coverage below 25). This result is consistent with previous observations that low coverage in sequences causes “breaks” in connectivity within the assembly graph (18, 19).

To evaluate partitioning and its separation of species further, we introduced spiked, simulated reads from several *Escherichia coli* genomes into the HGMC dataset. First, simulated reads from a single genome (*E. coli* strain E24377A, NC\_009801.1 with a 2% substitution error and 10× coverage) were added to the HGMC dataset, and the resulting dataset, HGMC.Ecoli1, was normalized by coverage, partitioned, and assembled. Similar amounts of data reduction were observed after digital normalization and partitioning (Table 1). Among the resulting 81,154 partitioned sets of reads in the HGMC.Ecoli1 dataset, only 2,580 partitions (3.2%) contained reads from multiple genomes. In total, 424 partitions contained reads from the spiked *E. coli* genome (201 partitions contained only spiked reads), and, when assembled, the contigs aligned to 99.5% of the *E. coli* strain E24377A genome (4,957,067 of 4,979,619 bp).

Next, we introduced five closely related *E. coli* strains [97.3–98.7% average nucleotide identity (20)] into the original HGMC



**Fig. 1.** Summary of approaches for large-scale assembly of complex metagenomes presented in this study. Unprocessed (I), normalized (II), and partitioned assemblies (III) were evaluated and compared with the HGMC metagenome. These approaches were used toward the assembly of metagenomes.



**Table 1. Total number of reads in unfiltered, normalized, and partitioned datasets**

Dataset	Unfiltered reads (Mbp)	Normalized reads (Mbp)	Partitioned reads (Mbp)
HGMC	14,494,884 (1,136)	8,656,520 (636)	8,560,124 (631)
HGMC spike	14,992,845 (1,137)	8,189,928 (612)	8,094,475 (607)
HGMC multispike	17,010,607 (1,339)	9,037,142 (702)	8,930,840 (697)
Iowa corn	1,810,630,781 (140,750)	1,406,361,241 (91,043)	1,040,396,940 (77,603)
Iowa prairie	3,303,375,485 (256,610)	2,241,951,533 (144,962)	1,696,187,797 (125,105)

dataset. This dataset, referred to as HGMC.Ecoli5, was normalized, partitioned, and assembled, resulting in 81,425 partitions. Among these, 1,154 partitions (1.4%) contained reads associated with multiple genomes. Among the partitions that contained reads associated with a single genome, 658 partitions contained reads originating from one of the spiked *E. coli* strains. In partitions containing reads from more than one genome, 224 partitions contained reads from a spiked *E. coli* strain and one other reference genome (either from another spiked strain or from the HGMC dataset). Independently assembling the partitions containing reads originating from the spiked *E. coli* strains resulted in 6,076 contigs, all but three originating from a spiked *E. coli* genome. The remaining three contigs were more than 99% similar to HGMC reference genomes (NC\_000915.1, NC\_003112.2, and NC\_009614.1). The contigs associated with the five *E. coli* strains aligned to more than 98% of each of the five genomes. Many of these contigs contained similarities to reads originating from multiple genomes found in the HGMC, and more than half of the contigs (3,075) could be aligned to reads that originated from more than one spiked genome.

For comparison, the HGMC.Ecoli5 dataset also was assembled without using any filtering approaches (e.g., no digital normalization or partitioning). In comparing the unfiltered and filtered HGMC.Ecoli5 assemblies, we found that the fractions of contigs associated with multiple genomes were similar. The assembly of the unfiltered dataset resulted in a greater proportion of contigs (66% or 4,702 contigs vs. 51% or 3,075 normalized/partitioned contigs) associated with multiple genomes.

**Assembly of Two Soil Metagenomes.** We next applied digital normalization and partitioning approaches to the de novo assembly of two soil metagenomes. Unfiltered Iowa corn and prairie datasets (containing 1.8 billion and 3.3 billion reads, respectively) could not be assembled by Velvet in 500 GB of RAM. A 75-million-reads subset of the Iowa corn dataset alone required 110 GB of memory, suggesting that assembly of the 3.3-billion-read dataset might need as much as 4 TB of RAM. Applying both normalization and partitioning approaches reduced the Iowa corn and prairie datasets to 1.4 billion and 2.2 billion reads, respectively, and after partitioning a total of 1.0 billion and 1.7 billion reads remained, respectively. These prefiltering approaches required 300 GB of RAM or less (Table 2). Notably, the large majority of k-mers in the soil metagenomes are relatively low abundance (Fig. 2), and consequently digital normalization did not remove as many reads in the soil metagenomes as in the mock dataset (Table 1).

Based on the HGMC dataset, we estimated that above a sequencing depth of five, the large majority of sequences that could be aligned to reference genomes are also assembled into contigs

greater than or equal to 300 bp (Fig. S2). Given the greater diversity expected in the soil metagenomes, we normalized these datasets to a sequencing depth of 10 (i.e., setting aside redundant reads within dataset above this coverage). After partitioning the filtered datasets, we identified a total 31,537,798 and 55,993,006 partitions (containing more than five reads) in the corn and prairie datasets, respectively. For assembly, we grouped partitions together into files containing a minimum of 10 million reads. Data reduction and partitioning were completed in less than 300 GB of RAM; once partitioned, each group of reads could be assembled in less than 14 GB and 4 h, readily enabling the evaluation of multiple assemblers and assembly parameters with practical computational resources.

The final assembly of the corn and prairie soil metagenomes resulted in a total of 1.9 million and 3.1 million contigs (minimum length of 300 bp), respectively, and a total assembly length of 912 million bp and 1.5 billion bp, respectively. To estimate abundance of assembled contigs and evaluate incorporation of reads, all quality-trimmed reads (including filtered reads) were aligned to assembled contigs. Overall, for the Iowa corn assembly, 8% of single reads and 10% of paired-end reads mapped to the assembly. Among paired-end reads, 95.5% of the reads aligned concordantly. In the Iowa prairie assembly, 10% of the single reads and 11% of the paired-end reads aligned to the assembled contigs, and 95.4% of the paired ends aligned concordantly (Table 4). Based on the alignment of sequencing reads to assembled contigs, we estimated the distribution of sequencing coverage in the resulting assemblies (Fig. 2). Overall, the coverage of each metagenome was low; 48% and 31% of total contigs in Iowa corn and prairie assemblies, respectively, had a read coverage less than 10.

Because the resulting assemblies are consensus representatives of the unassembled datasets, we also investigated the degree of variation (i.e., polymorphism) present among aligned reads to assembled contigs (*SI Materials and Methods*). For both the Iowa corn and prairie metagenomes, more than 99.9% of contigs contained base calls that were supported by a 95% consensus from mapped reads over 90% of their lengths, demonstrating an unexpectedly low polymorphism rate.

We annotated assembled contigs (greater than 300 bp) through the MG-RAST pipeline. This annotation resulted in 2,089,779 and 3,460,496 predicted protein coding regions in the corn and prairie metagenomes, respectively. The large majority of these regions, 61.8% in corn and 70.0% in prairie, had less than 60% similarity (over a minimum length of 15 aa) with any gene in the MG-RAST database M5NR (release 52). In total, 613,213 (29.3%) and 777,454 (22.5%) protein coding regions were assigned to an existing function. Many contigs were greater than 1 kbp, including 85,581 contigs in the corn metagenome (maximum length = 20,234) and 11,728 contigs in the prairie genome (maximum length = 2,579), and the distribution of lengths among assembled contigs was similar between sequences which could be assigned a function and those that could not (e.g., unknown sequences) (Figs. S3 and S4).

Annotations of the assembled corn and prairie soil metagenomes also were identified against the MG-RAST Kyoto Encyclopedia of Genes and Genomes Orthology (KEGG KO) database (Release 56). In total, 143,666 corn metagenome sequences and 164,318 prairie metagenome sequences matched sequences within the KO database with a minimum identity of

**Table 2. Computational resources (memory and time) required**

	Filter I: normalization, GB(h)	Filter II: partitioning, Gb(h)
HGMC	4(<2)	4(<2)
HGMC spike	4(<2)	4(<2)
HGMC multispike	4(<2)	4(<2)
Iowa corn	188(83)	234(120)
Iowa prairie	258(178)	287(310)

**Table 3. Assembly summary statistics for unfiltered, normalized, and normalized + partitioned datasets**

Dataset	No. contigs	Unfiltered length (Mbp)	Maximum contig (bp)	No. contigs	Normalized filtered length (Mbp)	Maximum contig (bp)	No. contigs	Partitioned length (Mbp)	Maximum contig (bp)	Assembler
HGMC	29,063	38	146,795	30,082	35	90,497	30,115	35	90,497	V
HGMC	24,300	36	86,445	—	—	—	27,475	36	96,041	M
HGMC	36,689	37	32,736	—	—	—	29,295	37	58,598	S
Iowa corn	—	—	—	—	—	—	1,862,962	912	20,234	V
Iowa corn	—	—	—	—	—	—	1,334,841	623	15,013	M
Iowa corn	—	—	—	—	—	—	1,542,436	675	15,075	S
Iowa prairie	—	—	—	—	—	—	3,120,263	1,510	9,397	V
Iowa prairie	—	—	—	—	—	—	2,102,163	998	7,206	M
Iowa prairie	—	—	—	—	—	—	2,599,767	1,145	5,423	S

M, MetalDBA assembler; S, SOAPdenovo assembler; V, Velvet assembler. Assemblies of Iowa corn and prairie metagenomes could not be completed on unfiltered or normalized-only datasets.

60%, a minimum length of 30 aa, and E-value <1e-10. The assembled contigs had significantly longer alignments to KEGG proteins than did unassembled reads ( $89 \pm 39$  aa vs.  $32 \pm 1$  aa) (Fig. S5). Among these, a total of 3,553 unique KO identifiers were identified (2,201 shared between corn and prairie metagenomes, 223 in corn alone, and 1,129 in prairie alone) and were found to represent broad metabolic functions (Fig. 3 and Fig. S6) involved in metabolism, genetic and environmental information processing, and cellular processes.

The shared presence of contigs without functional annotations in both the corn and prairie datasets also was evaluated. Assembled contigs that shared no homology to known sequences in the M5NR database were used as references for the complementing soil metagenome (e.g., corn assembly reference for prairie unassembled reads). Among these, a total of 34,436 contigs (31,058 and 3,416 corn and prairie contigs, respectively) were found to be shared between the two soil metagenomes (SI Materials and Methods).

## Discussion

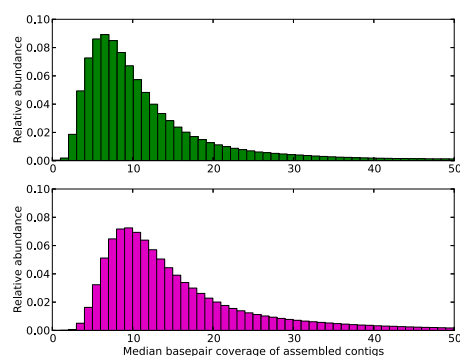
**Coverage-Based Filtering Approaches Reduce Datasets Without Information Loss.** Our described approach for scalable metagenomic assembly was effective in reducing the HGMC dataset size without significant loss of information. Although the diversity and sequencing depth represented by the HGMC dataset are extremely low as compared with most environmental metagenomes, it represents a simplified, unevenly sampled model for a metagenomic dataset that allows the evaluation of novel approaches through the availability of source genomes. Our approaches normalized the abundance of reads in the dataset to a specific sequencing coverage while reducing the dataset volume and removing errors introduced by extraneous reads. Furthermore, the partitioning approach subdivides large datasets into

biologically relevant subsets that can be assembled separately and by any assembler.

Based on our evaluations of a mock metagenome, we observed that the specific effects of filtering by digital normalization vary, depending on the conservation of genomic regions and abundance of genomes. In general, filtered (normalized and partitioned) assemblies were similar to or improved upon the assembly of the unprocessed dataset, suggesting that removing and subsetting these data do not result in substantial loss of information. For low-coverage genomes, the removal of erroneous sequences during normalization resulted in improved assemblies. The assembly of highly abundant genomes sharing conserved regions [such as the plasmids of the *Staphylococcus epidermidis* (NC\_005008.1, NC\_005007.1, and NC\_005003.1)] was negatively affected by normalization. The greater number of reads representing these sequences within the unfiltered mock metagenome likely enabled assemblers to extend the assembly of these sequences more effectively, and this advantage ultimately was observed as an increased recovery of these genomes in this assembly as compared with the normalized assembly. This result identifies a shortcoming of our approach for metagenomic assembly and, indeed, of most short-read assembly approaches, related to repetitive regions and/or polymorphisms. Although data reduction may cause some information loss, we exchanged this disadvantage for the ability to assemble previously intractable datasets. Our evaluation of the mock metagenome suggests that this information loss is minimal overall and that our approach results in a comparable assembly whose abundance estimations are slightly improved.

## Partitioning Separates Metagenomes into Tractable Subsets Representative of Species.

Metagenomes contain many distinct genomes that are largely disconnected from each other but that often share sequences as the result of sequence conservation or lateral transfer. Our prefiltering normalization approach removes both common multigenome elements and most artificial connectivity stemming from the sequencing process. The removal of these sequences does not significantly alter the recovery of HGMC reference genomes through de novo assembly, in which the resulting assemblies of unfiltered, normalized, and partitioned datasets were nearly identical. Further, for the mock metagenome, the large majority of partitions contained reads from a single reference genome, supporting our previous hypothesis that most connected subgraphs contain reads from distinct genomes (14). When an *E. coli* genome of 10× sequencing coverage was spiked into this dataset, it was divided into 424 partitions, likely because of the presence of introduced sequencing errors. Although fewer than half of these partitions ( $n = 201$ ) contained reads unique to the original reference genome, the combined assembly of each partition could recover nearly all of the original reference. When five similar *E. coli* genomes were mixed with the mock metagenome, we observed more partitions ( $n = 658$ ) containing *E. coli* sequences, one-third of which contained only *E. coli* sequences.



**Fig. 2.** Coverage (median base pair recovered) distribution of assembled contigs from the Iowa corn soil (Upper) and Iowa prairie soil (Lower) metagenomes.





of animal manure that potentially could enrich specific metabolic pathways with decreased diversity.

A key challenge facing future soil investigations is the lack of culturable representatives from soil and consequently the poor availability of reference genomes. This problem is highlighted by our observation that more than half of the assembled contigs were not similar to any sequence in the MG-RAST m5nr databases, suggesting that soil holds considerable unexplored taxonomic and functional novelty. These “unknown” sequences are broadly distributed in both length and abundance (Figs. S3 and S4) and represent the potential of gene and organism discovery. These sequences highlight the value of using de novo assemblies as reference datasets that are more representative of site-specific genes than are the publicly available references (where the average homology of assembled sequences against the SEED database was 68% over an average of 70 bp). For example, we identified 17 Mbp of unknown sequences in 34,436 contigs that were shared at relatively high abundance ( $C > 10$ ) by the corn and prairie soil metagenomes. These broadly present, novel sequences are targets for further investigations of proteins about which nothing is known. As increasing numbers of metagenomes become available, the co-occurrence of these assembled sequences with known genes and genomes will enable further characterization.

## Conclusions

We have presented two strategies that readily enable the assembly of very large environmental metagenomes by compressing and subdividing the data before assembly. The strategies are generic and broadly applicable to any metagenome. We demonstrate their effectiveness by first evaluating them on the assembly of a mock community metagenome and then applying them to two previously intractable soil metagenomes. Digital normalization scales the data size with community richness rather than diversity and is particularly effective for mixed-abundance communities. After digital normalization, partitioning enables the extraction of read subsets that belong to individual species. These read partitions are small enough that a variety of genomic-based analysis techniques can easily be applied to them individually, as evidenced by the application of multiple assemblers for our soil metagenomes with considerably reduced

computational resources. By acting as prefilters, digital normalization and partitioning let downstream assemblers focus on improving their performance on low-coverage or high-variability data without a strong consideration for computational resources. This ability should enable significant improvement of metagenome assembly techniques going forward and provide the critical references that will enable future investigations of soils and other complex environments. Importantly, our assembly results also demonstrate that 300 Gbp of read data are insufficient to cover even a small, localized soil sample deeply, confirming that considerably more data are needed to study the content of soil metagenomes comprehensively.

## Materials and Methods

Assemblies of the HGMC and soil metagenomes using various software were performed on (i) quality-filtered unassembled sequences and (ii) the same sequences filtered by digital normalization [HGMC coverage threshold ( $C$ ) = 20, soil coverage threshold ( $C$ ) = 10], removal of high-coverage sequences ( $C > 50$ ), and partitioning disconnected sets of reads. Coverage of assembled sequences or reference genomes was estimated through consensus alignment of raw sequences, and assembled contigs were compared with one another or reference genomes through BLASTn alignment (see *SI Materials and Methods* for specific thresholds). Annotation of assembled metagenomes and quality-filtered unassembled reads was performed through the MG-RAST and the M5NR (version 1) database and are available publicly (see *SI Materials and Methods* for accession numbers).

**ACKNOWLEDGMENTS.** We thank Krystle Chavarria and Regina Lamenella for help in extracting DNA from Great Prairie soil samples; Fan Yang for helpful comments on this paper; Eddy Rubin and Tijana Glavina del Rio at the Department of Energy Joint Genome Institute (DOE JGI); and John Johnson and Eric McDonald at the Michigan State University High Performance Computing Center. This project was supported by Agriculture and Food Research Initiative Competitive Grant 2010-65205-20361 from the US Department of Agriculture and by National Institute of Food and Agriculture and National Science Foundation Grant IOS-0923812 (both to C.T.B.). A.C.H. was supported by National Science Foundation Postdoctoral Fellowship Award 0905961 and the Great Lakes Bioenergy Research Center (Department of Energy BER DE-FC02-07ER64494). The work conducted by the DOE JGI is supported by the Office of Science of the US Department of Energy under Contract DE-AC02-05CH11231.

- Arumugam M, et al.; MetaHIT Consortium (2011) Enterotypes of the human gut microbiome. *Nature* 473(7346):174–180.
- Hess M, et al. (2011) Metagenomic discovery of biomass-degrading genes and genomes from cow rumen. *Science* 331(6016):463–467.
- Iverson V, et al. (2012) Untangling genomes from metagenomes: Revealing an uncultured class of marine Euryarchaeota. *Science* 335(6068):587–590.
- Mackelprang R, et al. (2011) Metagenomic analysis of a permafrost microbial community reveals a rapid response to thaw. *Nature* 480(7377):368–371.
- Qin J, et al.; MetaHIT Consortium (2010) A human gut microbial gene catalogue established by metagenomic sequencing. *Nature* 464(7285):59–65.
- Tringe SG, et al. (2005) Comparative metagenomics of microbial communities. *Science* 308(5721):554–557.
- Venter JC, et al. (2004) Environmental genome shotgun sequencing of the Sargasso Sea. *Science* 304(5667):66–74.
- Gans J, Wolinsky M, Dunbar J (2005) Computational improvements reveal great bacterial diversity and high metal toxicity in soil. *Science* 309(5739):1387–1390.
- Wommack KE, Bhavsar J, Ravel J (2008) Metagenomics: Read length matters. *Appl Environ Microbiol* 74(5):1453–1463.
- Loman NJ, et al. (2013) A culture-independent sequence-based metagenomics approach to the investigation of an outbreak of Shiga-toxicogenic *Escherichia coli* O104:H4. *JAMA* 309(14):1502–1510.
- Boisvert S, Raymond F, Godzaridis E, Laviolette F, Corbeil J (2012) Ray Meta: Scalable de novo metagenome assembly and profiling. *Genome Biol* 13(12):R122.
- Scholz MB, Lo C-C, Chain PSG (2012) Next generation sequencing and bioinformatic bottlenecks: The current state of metagenomic data analysis. *Curr Opin Biotechnol* 23(1):9–15.
- Brown CT, Howe A, Zhang Q, Pyrkosz AB, Brom TH (2012) A reference-free algorithm for computational normalization of shotgun sequencing data. *arXiv:1203.4802*. Accessed February 6, 2014.
- Pell J, et al. (2012) Scaling metagenome sequence assembly with probabilistic de Bruijn graphs. *Proc Natl Acad Sci USA* 109(33):13272–13277.
- Zerbino DR, Birney E (2008) Velvet: Algorithms for de novo short read assembly using de Bruijn graphs. *Genome Res* 18(5):821–829.
- Li R, et al. (2010) De novo assembly of human genomes with massively parallel short read sequencing. *Genome Res* 20(2):265–272.
- Peng Y, Leung HCM, Yiu SM, Chin FYL (2011) Meta-IDBA: A de novo assembler for metagenomic data. *Bioinformatics* 27(13):i94–i101.
- Chaisson MJ, Pevzner PA (2008) Short read fragment assembly of bacterial genomes. *Genome Res* 18(2):324–330.
- Pevzner PA, Tang H, Waterman MS (2001) An Eulerian path approach to DNA fragment assembly. *Proc Natl Acad Sci USA* 98(17):9748–9753.
- Goris J, et al. (2007) DNA-DNA hybridization values and their relationship to whole-genome sequence similarities. *Int J Syst Evol Microbiol* 57(Pt 1):81–91.
- Sharon I, et al. (2013) Time series community genomics analysis reveals rapid shifts in bacterial species, strains, and phage during infant gut colonization. *Genome Res* 23(1):111–120.

# Supporting Information

Howe et al. 10.1073/pnas.1402564111

## SI Materials and Methods

**Assembly Methods.** The Human Microbiome Project (HMP) mock community dataset and its available draft reference genomes were used to evaluate our approaches toward data reduction and partitioning for de novo metagenomic assembly. Reads of the mock community dataset initially were digitally normalized to a coverage threshold of 20 (as previously described in ref. 1), reducing the total number of reads from 14 to 11 million. Additionally, to remove possible sequencing artifacts associated with high-coverage sequences, highly abundant sequences (20-mers present at coverage greater than 50-fold) were filtered, and the dataset was normalized further to a coverage of 10, resulting in a total of 9 million reads (Fig. 3). Finally, the remaining reads were divided into disconnected sets of reads, resulting in a total of 85,818 partitions containing more than five reads (summarized in Table 1). The entire assembly pipeline for the mock community is described in detail in an IPython notebook available for download at <http://nbviewer.ipython.org/github/ngs-docs/ngs-notebooks/blob/master/ngs-71-hmp-partition.ipynb> and <http://nbviewer.ipython.org/github/ngs-docs/ngs-notebooks/blob/master/ngs-70-hmp-diginorm.ipynb>.

**Datasets.** In this study, we examined two large soil metagenomes generated from soils collected from Iowa corn and native prairie soils (6.4 and 3.2 g, respectively). Sequencing was performed at the Department of Energy Joint Genome Institute (Walnut Creek, CA). Soil metagenomes included both GAI and HiSeq Illumina paired-end libraries (76 × 2 and 114 × 2). Reads were quality trimmed at base pair locations where Phred scores indicated a score of 2. The total quality-trimmed reads in the Iowa corn and prairie datasets were 1.8 million and 3.3 million, respectively. Quality-trimmed unassembled reads used for Iowa corn and prairie assemblies are available on the MG-RAST metagenomics analysis server in project IDs 6368 (metagenomes 4539514.3–4539528.3) and 6377 (metagenomes 4539571.3–4539594.3). The direct urls for the Iowa corn metagenome and prairie unassembled reads: <http://metagenomics.anl.gov/linkin.cgi?project=6368> and <http://metagenomics.anl.gov/linkin.cgi?project=6377>. Untrimmed raw reads are also available at the National Center for Biotechnology Information Sequence Read Archive (SRA) as accession nos. SRX100357 and SRX099904–SRX099906 (Iowa Corn) and SRX100826, SRX099700, and SRX099701 (Iowa Prairie). We also include a human-gut mock community dataset (combined from SRA SRX055381 and SRX055380). For this mock community dataset, DNA from bacterial isolates originally recovered from within or on the human body were mixed together at staggered concentrations (over five orders of magnitude based on genomic DNA concentrations) and sequenced. The mock community dataset originally contained 14.5 million reads.

To evaluate our approach, we added simulated reads from either a single *Escherichia coli* (strain K-12 substrain DH10B) or five *E. coli* strains (K-12 substrains DH10B and E24377 and strain H7 substrains EC4115, UMN026, and SE15) into select metagenomes. We computationally generated 100 bp reads from each reference genome to a coverage of 10× and with a 2% error rate and subsequently randomly shuffled these reads.

**Estimation of Assembly Requirements for Soil Metagenomes.** Subsets of the Iowa corn metagenome were assembled with the Velvet assembler (v1.2.07) with the following parameters: `velveth K = 45, -short and velvetg -exp_cov auto -cov_cutoff auto, -scaf-`

folding no. The time and memory for each assembly was estimated up to a maximum of 150 h and 100 GB.

**Digital Normalization.** To reduce the dataset size, extraneous sequences for which sufficient coverage was available for assembly were removed through digital normalization. Digital normalization is described further in ref. 1. For the mock community dataset, digital normalization was performed with the following parameters:  $K = 20$ , coverage = 20, and Bloom filter size =  $1 \text{ GB} \times 4$ . For the Iowa corn metagenome, digital normalization parameters were as follows:  $K = 20$ , coverage = 10, and Bloom filter size =  $48 \text{ GB} \times 4$ . Similar parameters were used for the Iowa prairie metagenome, except that the Bloom filter size was  $60 \text{ GB} \times 4$ .

**Removal of High-Abundance Sequences.** To eliminate known sequencing artifacts in Illumina metagenomes, high-abundance sequences (coverage greater than 50) were removed using the count-min-sketch data structure used for digital normalization. For the relatively high-coverage mock community dataset, filtered reads were subsequently normalized to a coverage of 10 ( $K = 20$ , bloom filter size =  $1 \text{ GB} \times 4$ ).

**Partitioning and Assembly of Disconnected Reads.** The partitioning approach divides a larger dataset into groups of reads which share overlaps, dividing the dataset into subsets of connectivity. Disconnected partitions of the assembly de Bruijn graph (assembly graph) were separated by loading normalized filtered datasets into a probabilistic representation of the assembly graph as described in ref. 2. Partitions containing fewer than five reads were discarded. Each partition subsequently was assembled using the Velvet assembler with the setting as described above, with the exception that the  $K = 35$ –59 and shortPaired setting was used for paired-end reads. The resulting contigs greater than 300 bp from multiple- $K$  assemblies were dereplicated with CD-HIT [99% similarity (3)] and merged with Minimus2 (4). Additionally, partitions also were assembled with Meta-IDBA (5) and SOAPdenovo (6). Meta-IDBA assembly parameters were as follows: `idba -mink 25 -maxk 50 -minCount 0`. SOAPdenovo assembly parameters were as follows: `SOAPdenovo-31mer all -K 31 -p 8, asm_flags = 1`, and single and paired reads were separated and used for assembly.

**Comparing Coverage of Reference Genomes by Reads.** Reads in the HMP mock unfiltered and filtered datasets were mapped back to originating genomes using default settings in Bowtie2 (7). When reads could be mapped back to multiple genomes, a single genome was selected randomly to be identified with each read. Sequencing coverage was estimated for the whole genome as the median base pair coverage for all base pairs in the reference genome.

**Read Coverage by Assemblies.** All quality-trimmed reads, including those set aside by normalization, for Iowa corn and prairie were aligned with assembled contigs (length greater than 300 bp) using default parameters in Bowtie2 (7). Paired-end reads were evaluated according to concordance with paired-end library preparation (i.e., paired-end reads on opposite DNA strands) and the alignment of both pairs of reads to an assembled contig. The base pair coverage of each contig was estimated with the median base pair coverage of all reads across the length of the contig. Additionally, for each position in a contig (with the exception of the external 100 bp on each end), the percentage of the mapped consensus base pairs was calculated. The fraction of



positions with greater than 95% base consensus was calculated to estimate the presence of polymorphisms within the assembled contig.

Read coverage also was used to compare sequences without known annotations between the corn and prairie soil metagenomes. Contigs shared with another metagenome were characterized by having a minimum total of 10 bp median coverage among all available samples as well as being present in all samples. All identified shared contigs also were dereplicated with CD-HIT (99% similarity), and the best representative was chosen.

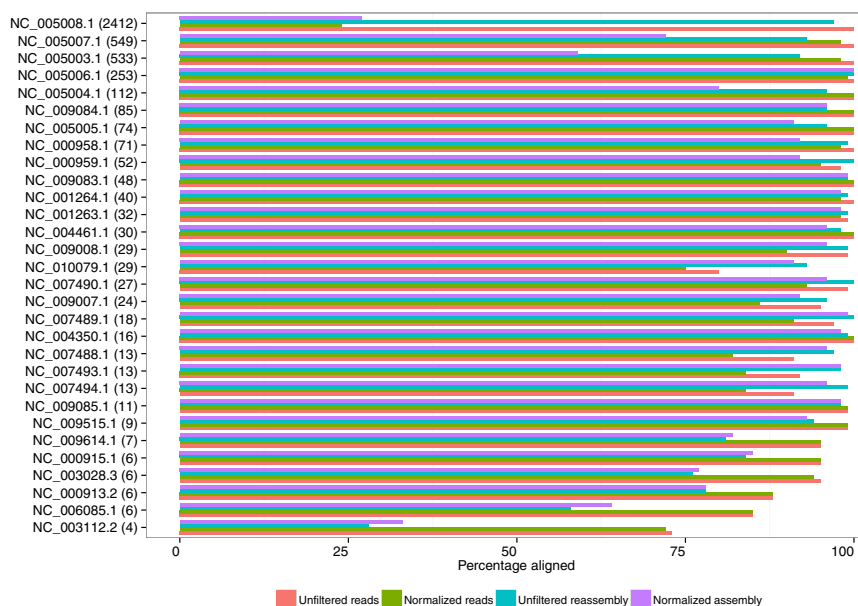
**Annotation of Assemblies.** Assembled contigs and their corresponding median base pair coverage for the Iowa corn and prairie metagenomes were uploaded into the MG-RAST annotation pipeline (8) and are available on MG-RAST as 4504979.3 (Iowa corn) and 4504798.3 (Iowa prairie). The resulting MG-RAST blast annotations were compared with the M5NR database using a maximum e-value of  $1e-5$ , a minimum identity of 60%, and a minimum alignment length of 15 aa unless otherwise noted. The annotated metagenome for Iowa corn can be found at [metagenomics.anl.gov/linkin.cgi?metagenome=4504797.3](http://metagenomics.anl.gov/linkin.cgi?metagenome=4504797.3) and for Iowa prairie at [metagenomics.anl.gov/linkin.cgi?metagenome=4504798.3](http://metagenomics.anl.gov/linkin.cgi?metagenome=4504798.3). The MG-RAST Kyoto Encyclopedia of Genes and Genomes Orthology (KEGG KO) ID numbers were placed on metabolic pathways available through [www.genome.jp/kegg/tool/map\\_pathway2.html](http://www.genome.jp/kegg/tool/map_pathway2.html).

**Comparing Assemblies.** Resulting assemblies (contigs greater than 300 bp) were compared using the total number of contigs, assembly

length, and maximum contig size for each assembly. Assemblies also were aligned to each other using BLASTn, and the resulting coverage of each assembly was calculated. For the mock community, the resulting assemblies also were aligned to sequenced draft genomes of the original isolates and, if applicable, to spiked reference genomes. Abundances of assembled contigs and reference genomes were estimated by mapping raw reads with Bowtie (allowing up to two mismatches for a match). The median base pair coverage was used to estimate abundances. Associated assembled contigs (greater than 300 bp) from the unfiltered and filtered (digital normalized) assemblies were identified using a BLASTn alignment (requiring an E-value cutoff of  $1e-5$ ). Contigs were associated with reference genomes through an identical alignment approach.

**Statistical Comparison of Assemblies.** The reference-based abundance (from reads mapped to reference genomes) and assembly-based abundance (from reads mapped to contigs) of genomes were compared. Using a one-directional, paired  $t$  test of squared deviations, the abundance estimates of the unfiltered and filtered assemblies were compared. The mean and SD of the abundances of unfiltered contigs, filtered contigs, and reference genes were  $6.8 \pm 7.1$ ,  $8.1 \pm 7.7$ , and  $7.8 \pm 5.2$ , respectively. We expected the filtered assembly to have increased accuracy because of a reduction of errors (e.g., normalization and high-abundance filtering) and used a one-sided  $t$  test which indicated that abundance estimations from the filtered assembly were significantly closer to predicted abundances from reference genomes ( $n = 28,652$ ,  $P = 0.032$ ).

1. Brown CT, Howe A, Zhang Q, Pyrkosz AB, Brom TH (2012) A reference-free algorithm for computational normalization of shotgun sequencing data. *arXiv:1203.4802*. Accessed February 6, 2014.
2. Pell J, et al. (2012) Scaling metagenome sequence assembly with probabilistic de Bruijn graphs. *Proc Natl Acad Sci USA* 109(33):13272–13277.
3. Fu L, Niu B, Zhu Z, Wu S, Li W (2012) CD-HIT: Accelerated for clustering the next-generation sequencing data. *Bioinformatics* 28(23):3150–3152.
4. Sommer DD, Delcher AL, Salzberg SL, Pop M (2007) Minimus: A fast, lightweight genome assembler. *BMC Bioinformatics* 8:64.
5. Peng Y, Leung HCM, Yiu SM, Chin FYL (2011) Meta-IDBA: A de novo assembler for metagenomic data. *Bioinformatics* 27(13):i94–i101.
6. Li R, et al. (2010) De novo assembly of human genomes with massively parallel short read sequencing. *Genome Res* 20(2):265–272.
7. Langmead B, Salzberg SL (2012) Fast gapped-read alignment with Bowtie 2. *Nat Methods* 9(4):357–359.
8. Meyer F, et al. (2008) The metagenomics RAST server - a public resource for the automatic phylogenetic and functional analysis of metagenomes. *BMC Bioinformatics* 9:386.



**Fig. S1.** Coverage of reference genomes by sequences in unfiltered and normalized filtered unassembled reads and assembled contigs. Estimated actual coverage is shown in parentheses.

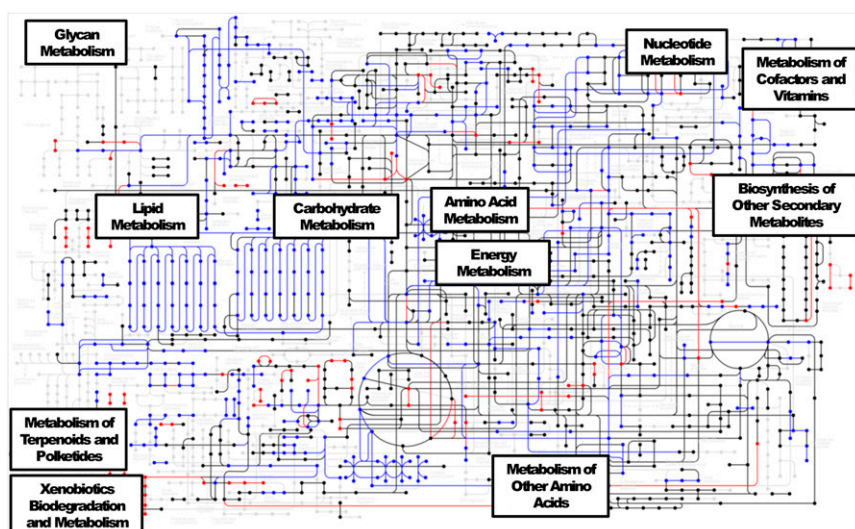
A scatter plot showing the relationship between the rank order of contigs (x-axis, 0 to 500) and their length in base pairs (y-axis, 0 to 20,000). The plot compares two sets of annotations: 'Known annotation' (red dots) and 'Unknown' (cyan dots). Each dot's size represents its  $\text{Log}(\text{Coverage})$ , with sizes corresponding to values 1.0, 1.5, 2.0, and 2.5. Both distributions show a rapid decrease in length as the rank order increases, with the 'Unknown' set generally having higher lengths than the 'Known' set for the same rank order.

3 of 6

Minimum Alignment Length	Unassembled Total Unique Hits	Assembled Total Unique Hits
10	129,205	79,574
20	130,016	79,574
30	68,037	79,477
40	0	77,148
50	0	68,168

4 of 6





**Fig. S6.** KEGG pathways sharing similarity with assembled soil metagenome sequences (black = combined corn and prairie, blue = prairie only, red = corn only).

Table S1. Human and Molecular Genome Center (HGMC) dataset reference genomes' estimated sequencing depth (coverage of reads expressed as median bp), number of partitions, total length, and coverage of reference genomes by unfiltered reads (UF cov), normalized filtered reads (F cov), unfiltered assembled contigs (UFA cov), and normalized filtered assembled contigs (FA cov)

Reference genome	Coverage	No. partitions	Length, bp	UF cov, bp	F cov, bp	UFA cov, %	FA cov, %
NC_005008.1	2,412	9	4,439	4,439	1,058	100	28
NC_005007.1	549	16	4,679	4,679	4,585	100	77
NC_005003.1	533	21	6,585	6,585	6,441	100	64%
NC_005006.1	253	2	8,007	8,004	7,953	100	100
NC_005004.1	112	52	24,365	24,358	24,291	100	83
NC_009084.1	85	3	11,302	11,295	11,270	100	100
NC_005005.1	74	12	17,261	17,202	17,180	100	100
NC_000958.1	71	73	177,466	177,261	174,614	100	95
NC_000959.1	52	37	45,704	44,974	43,557	100	92
NC_009083.1	48	2	13,408	13,405	13,383	100	100
NC_001264.1	40	63	412,348	410,970	403,553	100	99
NC_001263.1	32	546	2,648,638	2,634,512	2,589,566	100	99
NC_004461.1	30	476	2,499,279	2,498,081	2,492,248	100	98
NC_009008.1	29	14	37,100	36,585	33,250	94	96
NC_010079.1	29	442	2,872,915	2,298,758	2,157,196	100	92
NC_007490.1	27	27	100,828	99,385	93,550	100	96
NC_009007.1	24	92	114,045	108,526	97,860	100	96
NC_007489.1	18	12	105,284	102,212	96,169	100	99
NC_004350.1	16	131	2,030,921	2,029,376	2,025,544	100	99
NC_007488.1	13	30	114,178	103,351	93,637	100	99
NC_007493.1	13	628	3,188,609	2,919,441	2,681,855	100	99
NC_007494.1	13	262	943,016	862,781	788,626	100	98
NC_009085.1	11	683	3,976,747	3,939,190	3,936,208	99	99
NC_009515.1	9	552	1,853,160	1,828,231	1,826,639	99	98
NC_009614.1	7	7,751	5,163,189	4,899,622	4,896,808	81	82
NC_000915.1	6	2,888	1,667,867	1,581,502	1,581,024	78	79
NC_003028.3	6	4,123	2,160,842	2,047,832	2,037,347	78	78
NC_000913.2	6	5,913	4,639,675	4,080,605	4,074,119	84	85
NC_006085.1	6	6,459	2,560,265	2,169,547	2,169,056	59	64
NC_003112.2	4	9,269	2,272,360	1,655,023	1,626,301	28	33

**Table S2. Assembly comparisons of HGMC unfiltered (UF) and normalized filtered (NF) or filtered/partitioned (FP) HGMC datasets using different assemblers**

Assembly comparison	Similarity, %	RG coverage, %	Assembler
UF vs. NF	95	43.3/44.5	Velvet
UF vs. FP	95	43.3/44.4	Velvet
UF vs. FP	93	46.5/45.4	Meta-IDBA
UF vs. FP	98	46.2/46.4	SOAPdenovo

Assembly content similarity is based on the fraction of alignment of assemblies; similarly, the coverage of reference genomes (RG) is based on the alignment of assembled contigs to reference genomes.

**Table S3. Longest contigs in the corn and prairie soil metagenome with similarity to the RefSeq database**

Contig ID	Length, bp	Coverage	Function and organism
iowa-corn-3-pass.4582796.20234*	20,234	138	Probable poly(beta-D-mannuronate) O-acetylase (EC 2.3.1.-) <i>Cyanothece</i> sp PCC 7425
iowa-corn-3-pass.4606542.17507	17,507	30	Hypothetical protein <i>Pseudomonas</i> phage PaP2
iowa-corn-3-pass.4578484.16126	16,126	62	Hypothetical protein <i>Pseudomonas</i> phage PaP2
iowa-corn-3-pass.4583611.12814	12,814	18	Replication factor C large subunit <i>Candidatus Methanoregula boonei</i> 6A8
iowa-corn-3-pass.4594771.12496	12,496	35	Carbonic anhydrase <i>Psychromonas ingrahamii</i> 37
iowa-corn-3-pass.4596152.11816	11,816	25	Ribonuclease Z <i>Thermococcus kodakarensis</i> KOD1
iowa-corn-3-pass.4616349.11691	11,691	25	Precorrin-3B C17-methyltransferase <i>Nitrosopumilus maritimus</i> SCM1
iowa-corn-3-pass.4592007.10823	10,823	22	Hydroxymethylglutaryl-CoA reductase, degradative <i>Staphylothermus marinus</i> F1
iowa-corn-3-pass.4589906.10747	10,747	23	RdgB/HAM1 family noncanonical purine NTP pyrophosphatase <i>Roseiflexus castenholzii</i> DSM 13941
iowa-corn-3-pass.4559414.9998	9,998	19	Heparan N-sulfatase <i>Blastopirellula marina</i> DSM 3645
iowa-prairie-3-pass.6326293.1650	1,650	11	PAS/PAC sensor signal transduction histidine kinase <i>Marinobacter aquaeolei</i> VT8
iowa-prairie-3-pass.6215171.1615*	1,615	9	Beta-mannosidase (EC 3.2.1.25) <i>Acidobacteria bacterium</i> Ellin345
iowa-prairie-3-pass.6327344.1595	1,595	12	Glycosyl transferase, group 2 family protein <i>Bacillus cereus</i> ATCC 10987
iowa-prairie-3-pass.6326016.1586	1,586	10	Hypothetical protein <i>Dechloromonas aromatica</i> RCB
iowa-prairie-3-pass.6155396.1555	1,555	9	Peptidase S33, proline iminopeptidase 1 <i>Pseudomonas syringae</i> pv. <i>syringae</i> B728a
iowa-prairie-3-pass.6285889.1505	1,505	11	AMP-dependent synthetase and ligase <i>Chloroflexus</i> sp. Y-400-fl
iowa-prairie-3-pass.6293899.1502	1,502	13	Beta-lactamase <i>Flavobacterium johnsoniae</i> UW101
iowa-prairie-3-pass.5906709.1488	1,488	7	Peptidase M24 <i>Chitinophaga pinensis</i> DSM 2588
iowa-prairie-3-pass.6216765.1484	1,484	10	Amine oxidase <i>Candidatus Solibacter usitatus</i> Ellin6076
iowa-prairie-3-pass.6224149.1479	1,479	12	Phosphoribosylglycinamide formyltransferase <i>Bacteroides</i> sp. 2_1_16

\*Indicates SEED database origin if no RefSeq organismal annotation available.