

{Correspondence}

## Improving the usability and comprehensiveness of microbial databases

Caitlin Loeffler<sup>1†</sup>, Aaron Karlsberg<sup>1†</sup>, Lana S. Martin<sup>1</sup>, Eleazar Eskin<sup>1</sup>, David Koslicki<sup>2,3,4</sup>, Serghei Mangul<sup>5\*</sup>

<sup>1</sup> Department of Computer Science, University of California Los Angeles, 404 Westwood Plaza, Los Angeles, CA 90095, USA

<sup>2</sup> School of Computer Science and Engineering, The Pennsylvania State University, 207 Electrical Engineering West University Park, PA 16802, USA

<sup>3</sup> Department of Biology, The Pennsylvania State University, 208 Curtin Rd, State College, PA 16801, USA

<sup>4</sup> The Huck Institutes of the Life Sciences, The Pennsylvania State University, 101 Huck Life Sciences Building, University Park, PA 16802, USA

<sup>5</sup> Department of Clinical Pharmacy, University of Southern California School of Pharmacy, 1985 Zonal Ave, Los Angeles, CA 90089, USA

<sup>†</sup> These authors contributed equally to this work.

\*Correspondence: [mangul@usc.edu](mailto:mangul@usc.edu); [cloeffler@ucla.edu](mailto:cloeffler@ucla.edu)

## Abstract

Metagenomics studies leverage genomic reference databases to generate discoveries in basic science and translational research. We analyze existing fungal and bacterial databases and discuss guidelines for the development of a master reference database that promises to improve the quality and quantity of omics research.

## **Main Text**

High-throughput sequencing has revolutionized microbiome research by enabling the detection of thousands of microbial genomes directly from their host environments<sup>1</sup>. This approach, known as metagenomics, is capable of capturing the complex interactions that take place between thousands of different microbial organisms in their natural habitats. Metagenomic methods rely on comparisons of a sampled genome to multiple reference genomes. Metagenomics is more expensive to perform than traditional, culture-based taxonomic identification techniques, but today's metagenomic methods can produce a more comprehensive reconstruction of microbial genomes<sup>2</sup>. Emerging technologies for identifying and analysing microbial genomes can provide valuable insights into the interactions between human microbiomes and medicines. However, the current ad hoc practice of storing reference genomes in multiple, disparate reference databases challenges the accuracy and comprehensiveness of future microbial metagenomics studies.

Metagenomic studies isolate DNA found in a sample of various environments, compare the sampled genomes (represented as a set of reads) to verified reference genomes, and identify the organism from which the reads originated. Ideally, a metagenomic study uses a reference database that contains all known genomic

references. Today's researcher can choose from many different genomic reference databases that contain verified reference genomes, but these databases lack a universal standard of specimen inclusion, data preparation, taxon labeling, and accessibility.

Several limitations in genomic sequencing present unique challenges to accurately assembling reference genomes and compile them into comprehensive databases. Notably, reference genomes can exist in various stages of completion. Typically, reads are assembled into larger sequences which represent complete or fragmented microbial genomes. Fragmented assemblies are usually represented as a set of contigs, which are typically contiguous DNA fragments corresponding to unlocalized segments of microbial genomes. Given sufficient data, contigs can be further assembled into scaffolds that represent larger portions of individual chromosomes but gaps (consisting of a possibly unknown number of unknown nucleotides) can remain. Most reference genomes are in different stages of completeness, with portions of even the human genome remaining unknown (in particular, the centromere and telomere regions).

In addition, the location of possibly incomplete reference genomes on taxonomic or phylogenetic trees can be contentious. Metagenomics researchers must take into account discrepancies in the types of taxa included in each reference genome database, as well as differences in how the genomes are constructed, identified, and made available for distribution.

The future of metagenomics research would benefit from a standardized, comprehensive approach to reference genome database development. To begin

assembling a set of recommendations for reference genome database construction, we assessed the concordance and usability of available reference databases for microbial genomics. Our study considered the concordance of microbial species and genera across four fungal reference databases (Ensembl<sup>3</sup>, RefSeq<sup>4</sup>, JGI's 1000 fungal genomes project (JGI 1K)<sup>5</sup>, and FungiDB<sup>6</sup>) and three bacterial reference databases (Ensembl<sup>3</sup>, RefSeq<sup>4</sup>, and PATRIC<sup>7</sup>). We compared the microbial taxa in each of the databases using NCBI's universal taxonomic identifiers (hereafter referred to as taxIDs) at the ranks of species and genus. Strains were not included in this analysis as studied databases contained multiple instances where a reference was counted as a strain in one database yet was labeled an isolate in NCBI; in such cases, the reference was not yet assigned a strain-level NCBI taxID. This discrepancy made comparison of strain comprehensiveness among databases impossible to calculate and demonstrates the importance of developing a standardized taxonomic naming system to be shared between databases<sup>8</sup>.

Our comparison of four major fungal and three major bacterial genome databases reveals substantial discrepancies across databases in the presence of microbial references at taxonomic levels below the family rank. In other words, a researcher's selection of one particular reference database could substantially impact the number and types of unique microbial taxa identified in a study.

Calculating the coverage of each fungal reference genome database shows that a researcher using the largest—and most comprehensive— fungal reference database would only find identification for 89% of the total 786 fungal genera covered by all four databases (Figure 1c), which is 80% of the possible 1405 fungal species (Figure

1a). Similarly, calculating the coverage of each bacterial reference genome database shows that a researcher using the largest—and most comprehensive—reference database would find identification for 94% of the total 3371 bacterial genera covered by all three databases (Figure 1d), which is 95% of the possible 42,337 bacterial species (Figure 1b).

Only a relatively small percentage of species are represented as complete genomes; calculating the percentage of fungal species per reference database reveals that 16% of species are represented as complete fungal genomes in Ensembl, 2% in JGI 1K, 14% in RefSeq, and 13% in FungiDB. Conversely, our study shows that the percentage of species represented as contigs are relatively high: 81% in RefSeq, 98% in JGI 1K, 80% in Ensembl, and 81% in FungiDB. Remaining genomes are comprised of contigs or a mixture of chromosomes and contigs (Figure 2a). In addition, we found that complete reference genomes for fungi taxa were not consistently present in studied fungal reference databases. In total, there are 53 unique species represented across the four fungal databases that are complete genomes. Of these, only 13% are represented in all four databases (Figure S3).

We found similar results for bacterial species in the bacterial genome reference databases. Only 11% of bacterial references are represented as complete bacterial genomes in Ensembl, 10% in RefSeq, and 3% in PATRIC. The majority of references are represented as contigs in Ensembl (89%), Refeq (90%), and PATRIC (97%). All three bacterial genome reference databases have <1% of references containing a mix of contigs and chromosomes (Figure 2b).

Of the 80 - 90% of the references in each database represented as fragmented genomes, we considered the length distributions of the sequences provided. The length distributions for contigs are relatively similar across all four fungal databases (Figure 2c). The length distributions for contigs are relatively similar across the three bacterial databases we studied (Figure 2d). The mean contig length is shorter in bacterial reference databases than in fungal reference databases.

The completeness of a reference database is always subject to limitations imposed by the project's funding or scope. As one example of the latter, the JGI 1K database contains many novel and previously unpublished genomes. The introductory text of the JGI database indicates that, for this reason, it is not designed to be used in metagenomics studies<sup>5</sup>. However, such a large database of novel references may be a top choice for metagenomics researchers who want to learn as much as they can about their samples. Of the four fungal reference databases analyzed in this study, JGI 1K is the largest, covering 89% of fungal genera and 80% of fungal species. Ensembl, the second largest of the four databases, only covers 45% of fungal species and 41% of fungal genera.

In some cases, a more complete database may hinder analytical methods. Due to limitations in metagenomic analysis pipelines, reference databases containing species whose genomes are remarkably similar often prevent identification at the species level<sup>9</sup>.

Even taking these limitations into account, researchers would benefit from a universal approach to constructing comprehensive microbial genomic reference databases.

Since the ideal reference database containing all the reference genomes for all known samples does not yet exist, researchers are potentially failing to identify key organisms within their samples. The first consideration of a master reference database would be developing a standardized approach to assembling and presenting data from existing reference databases. A systematic approach to constructing reference databases, when adopted by the scientific community, would help improve microbial coverage in newly developed metagenomic analysis tools.

One approach to developing a comprehensive database of complete genomic references is to combine all existing reference databases into one master set—a complex, time-consuming task. With this approach, references unique to one database could simply be added to a master set. However, a reference that is found in more than one database presents several problems. Multiple references may be assigned the same taxID, yet these references may contain differing genomic information. For example, references comprised of contigs could cover different segments of a given gene. Selecting both unedited contig-based reference genomes would unnecessarily extend the run time of a comparison algorithm utilizing the master set. On the other hand, eliminating one reference would ignore entire segments of the genome represented in the discarded contigs. In such cases, the database developer needs a consistent method for selecting one of the references to include in the master set.

An alternative approach would be to develop an open source computational method that continuously merges any number of disjointed microbial reference databases as new sequences become available. The sequencing and storing of microbial species in multiple repositories present an opportunity to improve sequence quality through an

approach based on alignment and consensus. An open source format would encourage computational developers to contribute to the reference database by engineering support for the integration of other, lesser known, reference sequence repositories.

Another potential strategy is to eliminate discrepancies between databases. This will require the development of a communication protocol that allows databases to share information and complement each other in real time. Such a communication protocol could eventually enable an assembly of a comprehensive ‘virtual’ database, which essentially represent a consensus across databases. Several technical issues may pose difficulties in implementing such an approach. For example, the proposed approach needs to be capable of resolving the conflicts between the databases, such as when references are represented by different contigs across databases.

We would also like to mention that, just like reference databases for genetic data, the reference databases for taxonomies also have restricted overlap<sup>8</sup>. For this present study, we were able to use NCBI Universal Taxonomic IDs (taxIDs) to measure species and genus reference congruence across the databases since NCBI taxIDs were used by each database we studied. Hence, database discrepancies only existed due to presence or absence of organisms in the reference database, not due to taxonomic ambiguities. However, there exist many such universal taxonomic systems which may overlap very little and where there may not exist a mapping to convert from one taxonomy to another. Further, even though we were able to identify species and genus across databases by NCBI taxIDs, this did not extend to strains as NCBI does not universally assign taxonomic identifiers to strains. The master database for reference



genomes, will, therefore, also need to utilize a master database for taxonomy. For example, one possible master taxonomic database may be the OpenTree taxonomy<sup>8</sup>.

A second consideration of a master reference database is usability. Bioinformatics is an interdisciplinary field comprised of researchers with varied backgrounds—from computer science to biology. In order to maximize potential use by both skilled and novice computational users, this complete database would need an intuitive user interface.

The four fungal and three bacterial databases analysed in this study presented challenges to data access and manipulation. For example, the fungal JGI 1K asks the user to select the genomes of interest from a picture of the fungal tree of life, which can be unintuitive to many researchers. Adequate user support would also increase the usability of a comprehensive reference database; at the time of our study, Ensembl did not publish any contact information on their webpage.

Several reference databases highlight features would be helpful to implement in a master reference database. FungiDB's interface, which is more intuitive, simply asks the user to select data as though shopping online. To download all organisms, one only had to hover over "About FungiDB", click "Organisms" under "----- Data in FungiDB", click "add to basket". Once all the organisms are placed in the basket, it is possible to customize an annotation table containing download links for all references within the basket. While downloading data from NCBI RefSeq can be challenging, once the user knows to select "Assembly" in the dropdown menu on the home page and type "Fungi" into the search bar, the filtering process becomes more intuitive. Th

“shopping basket” method is not efficient for downloading bacterial references, however, as there are over 200 thousand references to handle. A better approach would be to allow the user to download references from the NCBI FTP site, which requires knowledge of the command line and may not be usable by researchers lacking a computational background.

A third consideration of a master reference database is maintenance support and archival stability. Maintaining a master reference sequence database would carry a substantial cost in terms of computational power and storage. An open source, continuous assembly approach would depend on support from an institution, governing body, or a global consortium.

The Pathosystems Resource Integration Center (PATRIC) online bacterial reference database can be used as a gold standard for database website design. In PATRIC, all genomes for the selected taxa are present, and the filtering is intuitive. One drawback to the PATRIC website is the current protocol for downloading genomes; the best way to transfer data between servers is to generate a list of genome\_ids in the command line for the genomes of interest, then recursively call “wget” on each genome. Any researcher not familiar with the command line needs to download the data directly from the PATRIC website; this method is not allowed for bulk downloads. A more efficient alternative method for bulk downloading reference data without using the command line would be to provide an option to utilize a data transfer service (such as Globus), which PATRIC does not currently use.

Our study indicates that the current approach to developing genomic reference databases for fungal and bacterial species are not meeting the needs of metagenomics research. As metagenomic data become increasingly high-resolution, researchers need tools that enable a similarly high precision of taxonomic identification of DNA derived from samples. We believe that a systematic approach to developing a centralized master reference database will increase coverage and dramatically improve the quality and quantity of -omics research.

### ***Declarations***

### **Acknowledgements**

Not Applicable

### **Funding**

Not Applicable

### **Availability of Data and Materials**

The data supporting the conclusions of this article, including the species and genera names, are available at: <https://github.com/Mangul-Lab-USC/db.microbiome>.

### **Competing interests**

The authors declare that they have no competing interests.

### **Author contributions**

C.L. and A.K. performed the analysis. C.L., L.M., D.K., and S.M. wrote the manuscript. S.M. conceived of the presented idea. E.E. supervised the project. All authors have read and agree to the content.

## References

### There should be no more than 10 references in a Comment

1. Sharpton, T. J. An introduction to the analysis of shotgun metagenomic data. *Front. Plant Sci.* **5**, 209 (2014).
2. Hilton, S. K. *et al.* Metataxonomic and Metagenomic Approaches vs. Culture-Based Techniques for Clinical Pathology. *Front. Microbiol.* **7**, 484 (2016).
3. Kersey, P. J. *et al.* Ensembl Genomes 2018: an integrated omics infrastructure for non-vertebrate species. *Nucleic Acids Res.* **46**, D802–D808 (2018).
4. O’Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res.* **44**, D733–45 (2016).
5. Nordberg, H. *et al.* The genome portal of the Department of Energy Joint Genome Institute: 2014 updates. *Nucleic Acids Res.* **42**, D26–31 (2014).
6. Basenko, E. *et al.* FungiDB: An Integrated Bioinformatic Resource for Fungi and Oomycetes. *Journal of Fungi* vol. 4 39 (2018).
7. Wattam, A. R. *et al.* Improvements to PATRIC, the all-bacterial Bioinformatics Database and Analysis Resource Center. *Nucleic Acids Res.* **45**, D535–D542 (2017).
8. Hinchliff, C. *et al.* Synthesis of phylogeny and taxonomy into a comprehensive tree of life. doi:10.1101/012260.
9. Nasko, D. J., Koren, S., Phillippy, A. M. & Treangen, T. J. RefSeq database growth influences the accuracy of k-mer-based lowest common ancestor species identification. *Genome Biol.* **19**, 165 (2018).

# **{Figure Legends}**

Figure 1.

Consensus of fungal and bacterial genome representation across multiple reference databases.

(a) In total, there are 1405 unique species represented across the four fungal databases. Of these, 48 species are represented in all four databases. There are a total of 175 species found where strictly three databases overlap and 189 species where strictly two databases overlap. A total of 993 unique fungal species cannot be found in any overlaps. (b) In total, there are 42337 unique species represented across the three bacterial databases. Of these, 6543 species are represented in all three databases, and 17506 total species are found where strictly two databases overlap. A total of 18288 unique bacterial species cannot be found in any overlaps. (c) In total, there are 786 unique genera represented across the four fungal databases. Of these, 29 genera are represented in all four databases. There are a total of 109 genera found where strictly three databases overlap and 142 genera where strictly two databases overlap. A total of 506 unique fungal genera cannot be found in any overlaps. (d) In total, there are 2214 unique genera represented across the three bacterial databases. Of these, 76 genera are represented in all three databases, and 1149 total genera are found where strictly two databases overlap. A total of 989 unique bacterial genera cannot be found in any overlaps.

Figure 2.

Fungal and bacterial genome composition across multiple reference databases.

(a) Percentage of references per fungal database available as complete genomes (yellow), fragmented genomes (i.e., set of contigs) (blue), and a mixture of full chromosomes and contigs (red). (b) Percentage of species per bacterial database available as complete genomes (yellow), fragmented genomes (i.e., set of contigs) (blue), and a mixture of full chromosomes and contigs (red). (c) Length distribution of the fungal genomes (contigs) across the databases. The contig mean lengths for each fungal database are 322 thousand bp (Ensembl), 513 thousand bp (RefSeq), 426 thousand bp (JGI 1K), and 548K (FungiDB). (d) Length distribution of the bacterial genomes (contigs) across the databases. The contig mean lengths for each bacterial database are 149 thousand bp (Ensembl), 119 thousand bp (RefSeq), and 107 thousand bp (PATRIC).

## Figures

Figure 1

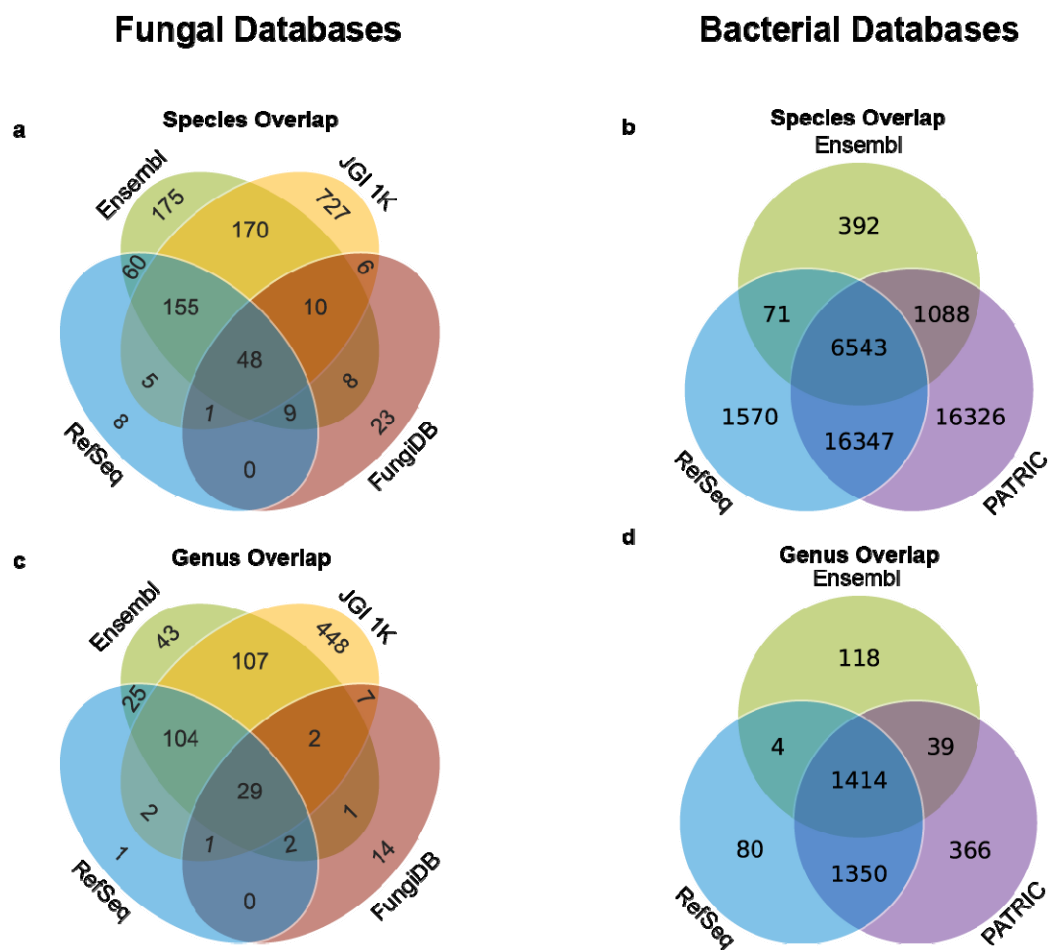
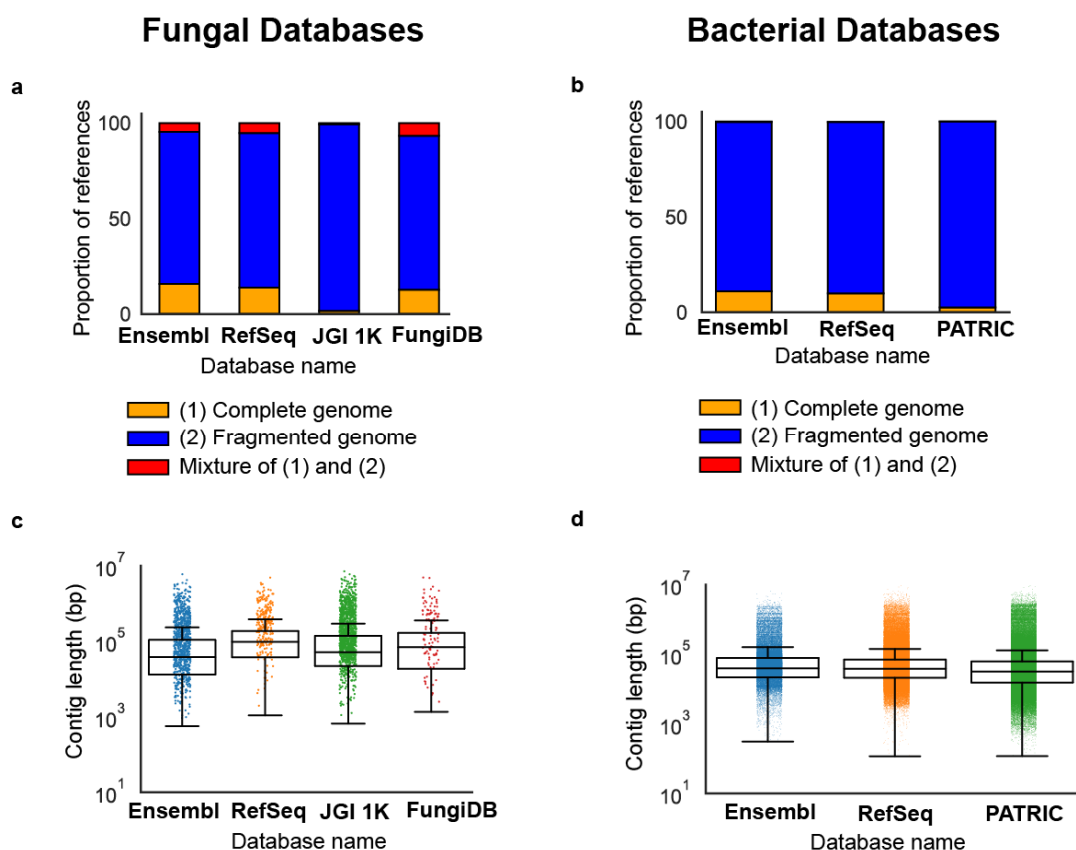


Figure 2



## Tables

Not Applicable

## Supplementary Materials

### Part I: Downloading the databases

Our study considered fungal species and genera across four reference databases:

- JGI 1000 Fungal Genomes Database (JGI 1K),  
<https://genome.jgi.doe.gov/programs/fungi/index.jsf>
- Ensembl, <http://fungi.ensembl.org/index.html>



- RefSeq, <https://www.ncbi.nlm.nih.gov/>
- FungiDB, <http://fungidb.org/fungidb/>

Our study also considered bacterial species and genera across three reference databases:

- Ensembl, <http://bacteria.ensembl.org/index.html>
- RefSeq, <https://www.ncbi.nlm.nih.gov/>
- PATRIC, <https://www.patricbrc.org/>

Each reference database had a different process for downloading reference genomes.

For fungi:

- **JGI 1K Fungal Genomes Database.** We locally downloaded the assembled masked fungal reference database, which appeared in the “download” section of the website as a zipped file. The unzipped file yielded 1265 directories. Each directory represented one species or, when strain information was available, one strain. The contents of each directory included a zipped FASTA file (in two directories there were two such files) that contained the genetic reference information. Plasmid and mitochondrial sequences were available in separate files.
- **Ensembl.** We called `wget` recursively on all FTP files in release 44 ending in `‘.dna.toplevel.fa.gz’`.
- **RefSeq.** We downloaded the FASTA files locally after filtering the NCBI assembly “Fungi” results for latest RefSeq references.
- **FungiDB.** Once all the fungal reference genomes were placed into the online basket, we downloaded the links corresponding to FASTA files. We then created a bash file that called “`wget`” on each link.

For bacteria:

- **Ensembl.** We ran a simple one line script accessing the FTP site
- **RefSeq.** We ran a recursive wget function that downloaded all the bacterial fasta files.
- **PATRIC.** The FTP site does not group genomes by lineage, and bulk downloads from the website itself is limited to sets of 10 thousand references at a time. A list of genome ids was generated in a text file which was then accessed by a wget loop to download each individual genome of interest.

Part II: Standardizing the taxonomy across the fungal and bacterial reference databases

In order to standardize the taxonomy across all four fungal and three bacterial reference databases, provided NCBI universal taxIDs were used in place of scientific names. TaxIDs are given at each taxonomic level and, therefore, can be ranked from Superkingdom to Strain levels. Only species- and genus-level taxIDs were used to quantify the consensus of fungal and bacterial genome representation across the databases. We did not analyze the consensus of strain-level taxIDs. We used the Ete3 module to assign a species-level taxID when the database provided a strain-level taxonomic identification, and to obtain a genus-level taxID for all reference genomes. We encountered multiple genomes that had not been assigned a genus and lacked a genus-level taxID; in the data, such reference genomes indicated ‘no rank’ where the genus-level taxID would have appeared. In such cases, we used the unranked taxID as the genus taxID.

As with the procedures for downloading reference genomes, the processes for obtaining corresponding taxIDs for each file were different for each of the databases.

- **JGI 1K Fungal Genomes Database.** We followed the six-step process for obtaining a Microsoft Excel document with the taxIDs. First, we created an advanced search that produced a “reports” button where the user can download the taxID information via an Excel spreadsheet. Second, we prepared a Python script to automatically match filenames with corresponding taxIDs.
- **Ensembl.** The FTP server (<ftp://ftp.ensemblgenomes.org/pub/fungi/release-44/>) offers a file named “species\_EnsemblFungi.txt,” a mapping file that shows which Ensembl files match to corresponding taxIDs. We ran a Python script that used the mapping file and a list of Ensembl file names to assign taxIDs.
- **RefSeq.** Similar to Ensembl, we used a mapping file to match taxIDs to accession numbers listed in the first header of each reference FASTA file.
- **FungiDB.** After downloading all the reference species, a csv file could be custom generated by clicking the “download” link, selecting “choose columns”, and checking the “NCBI taxon ID” box under “taxonomy”.

For the bacterial databases we isolated a list of taxIDs present in a given bacterial database and did not match taxIDs to filename.

Part III: Classify reference genomes as complete or fragmented

In order to determine the assembly level (e.g., scaffolds, contigs, fully assembled chromosomes) and identity of extra genetic material (e.g., mitochondrial and plasmid sequences) for each reference genome, we searched the headers of each reference FASTA file for predetermined patterns and words. We used the substrings “chr”, “complete”, and “NC\_” to identify sequences that had been marked as complete genomes. The key substrings “contig”, “scaffold”, “partial”, “supercont”, “unitig\_”, “NW\_”, and “NT\_” were used to identify sequences that had been marked as fragmented genomes.

When identifying the extra genetic material, we used the substrings “mitochondri” and “mt dna” to identify sequences that had been marked as mitochondrial references. Finally, we used the word “plasmid” to identify sequences that had been marked as plasmids.

#### Part IV: Compare the species and genera across the reference databases

In order to generate statistical data for cross-database reference comparison, we extracted individual sequence attributes from each reference FASTA file. We stored these attributes in a structured query language relational database management system (SQL RDBM). Attributes extracted from each fungal reference sequence included database name; species-level taxID; genus-level taxID; species name; genus name; a flag indicating reference composition (e.g., chromosomes, contigs or mixture of both); a flag indicating if the reference contains mitochondrial and plasmid DNAs.

For each reference, we also recorded the length of contigs and chromosomes.

Individual files could have more than one sequence classification depending on the contents of the DNA sequences. The data for sequence composition contained the

number of sequences for a given sequence classification that existed within each file.

We also stored within each file the average, minimum, and maximum sequence lengths for each sequence classification.

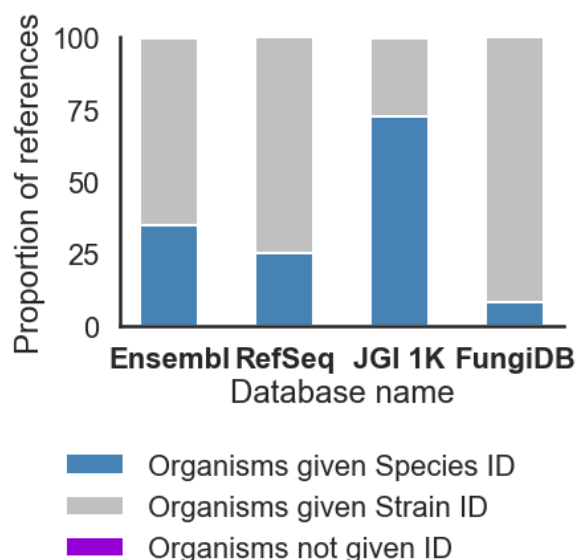


Figure S1. Proportion of fungal references that were identified by species-level or strain-level taxID, or were missing a taxID across three of the four databases. Species-level taxIDs were given to 35% of references in Ensembl, 25% in RefSeq, 73% in JGI 1K, and 8% in FungiDB. Strain-level taxIDs were given to 65% of references in Ensembl, 75% on RefSeq, 27% in JGI 1K, and 92% in FungiDB. TaxIDs were invalid or missing from < 0.5% of references in all four databases.

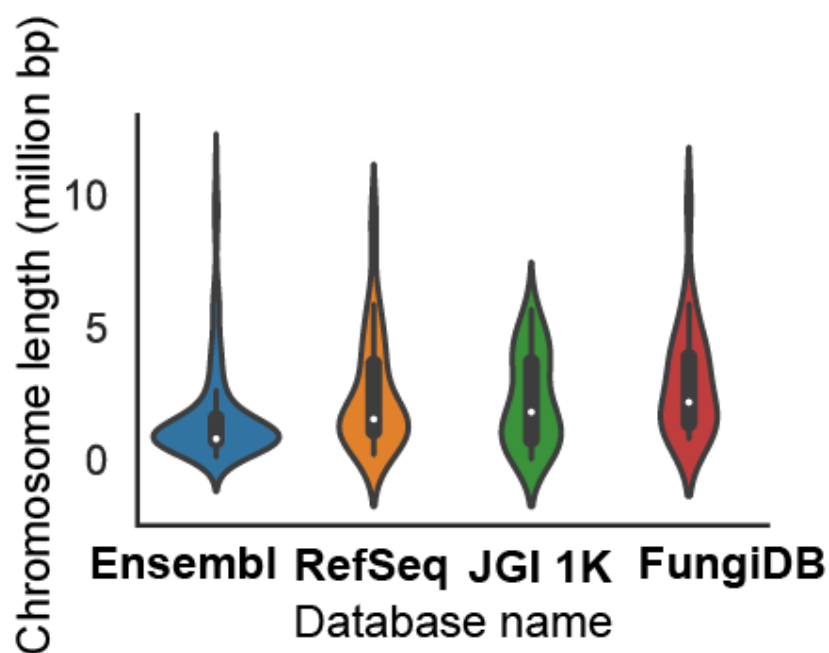


Figure S2. Length distribution of the fungal genomes (represented as chromosomes) across the databases. The mean lengths of chromosomes for each database is 1.7 million bp (Ensembl), 2.5 million bp (RefSeq), 2.2 million bp (JGI 1K), and 2.9 million bp (FungiDB).

## Overlap of complete genomes

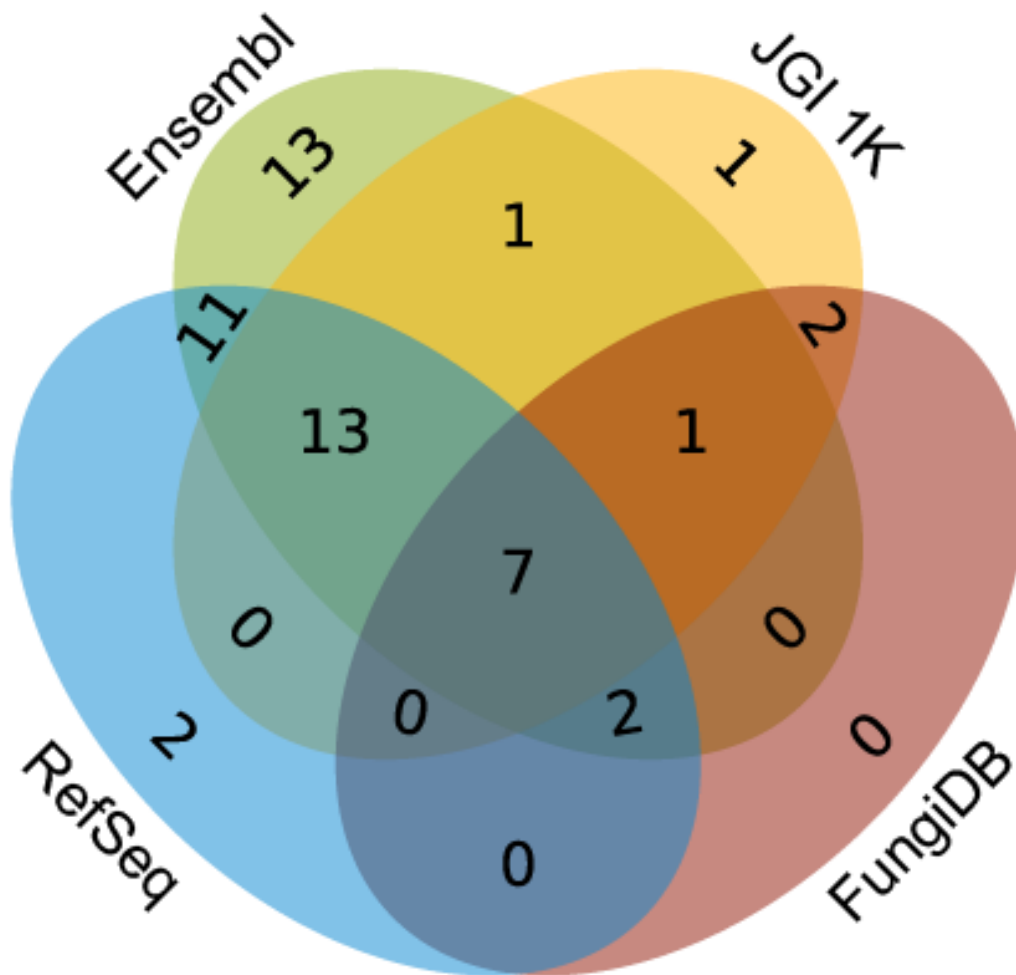


Figure S3. Overlap of species level references containing only complete chromosomes. In total, 53 unique species references contain only complete chromosomes represented across the four fungal databases. Of these, seven species are represented in all four databases. There are a total of 16 species found where strictly three databases overlap and 13 species where strictly two databases overlap. A total of 17 unique fungal species cannot be found in any overlaps.

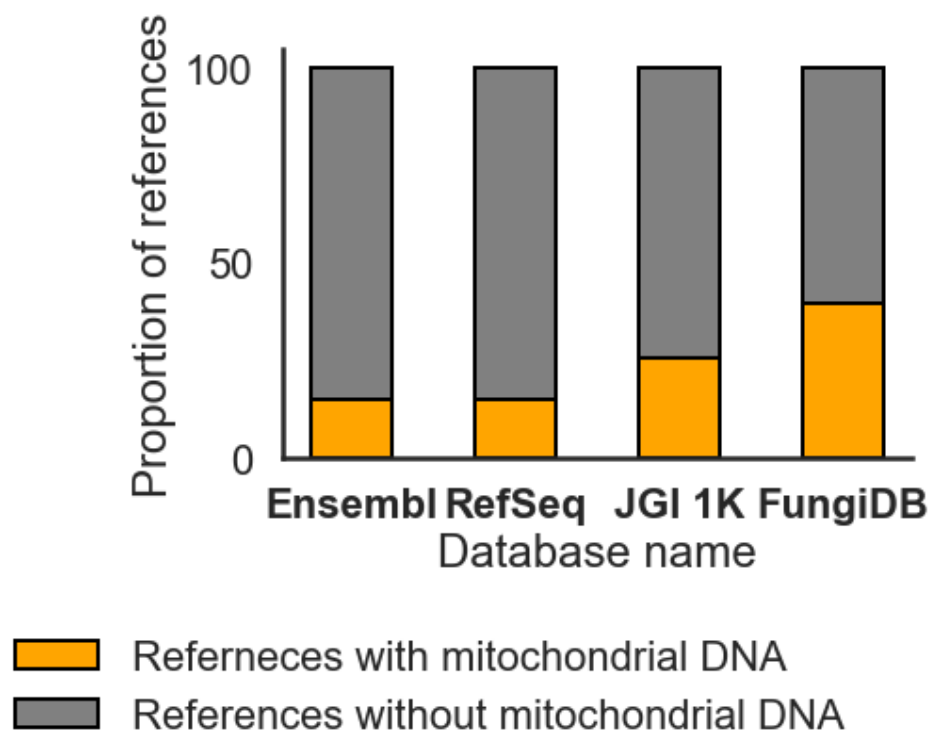


Figure S4. Percentage of references per database containing mitochondrial DNA (mtDNA) (orange). The percent of references that contained mitochondrial sequences are 15% (Ensembl), 15% (RefSeq), 25% (JGI 1K), and 40% (FungiDB).



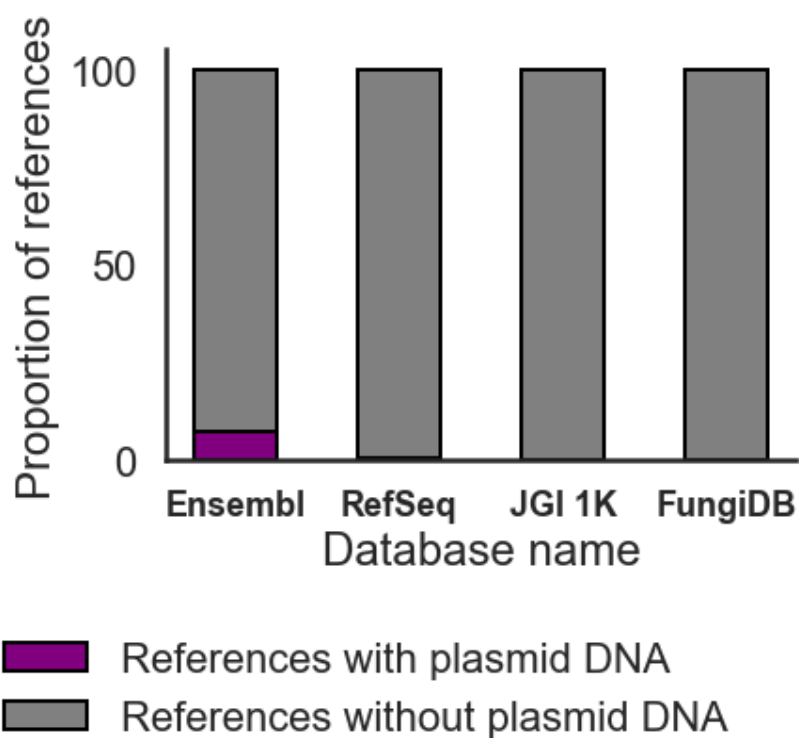
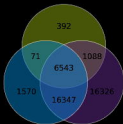


Figure S5. Percentage of references per database containing plasmid DNA (violet).

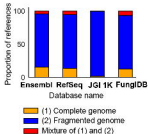
These percentages are 7% (Ensembl), 0.36% (RefSeq), 0% (JGI 1K), and 0%

FungiDB.



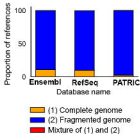
## Fungal Databases

a

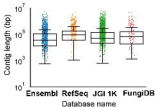


## Bacterial Databases

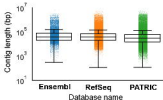
b

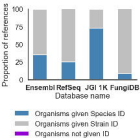


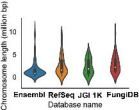
c



d

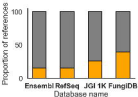






# Overlap of complete genomes

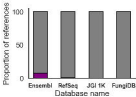




References with mitochondrial DNA



References without mitochondrial DNA



References with plasmid DNA



References without plasmid DNA