


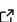
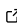
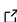
A Data Carpentry- Style Metagenomics Workshop

Claudia Ziri3n-Mart3nez¹, Diego Garfias-Gallegos¹, Tania Vanessa Arellano-Fernandez^{3,6}, Aar3n Espinosa-Jaime⁶, Edder D Bustos-D3az², Jos3 Abel Lovaco-Flores^{4,6}, Luis Gerardo Tejero-G3mez⁵, J Abraham Avelar-Rivas^{3,4}, and Nelly S3lem-Mojica ⁵

1 Laboratorio de Gen3mica Ecol3gica y Evolutiva, Langebio, Cinvestav, M3xico. 2 Laboratorio de Evoluci3n de la Diversidad Metab3lica, Langebio, Cinvestav, M3xico. 3 Laboratorio de Sistemas Gen3ticos, Langebio, Cinvestav, M3xico. 4 BetterLab - C3. Irapuato, M3xico. 5 Centro de Ciencias Matem3ticas. UNAM, M3xico. 6 Escuela Nacional de Estudios Superiores, Unidad Le3n, UNAM, M3xico.

DOI: [10.21105/jose.00209](https://doi.org/10.21105/jose.00209)

Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Submitted: 12 April 2023

Published: 01 September 2023

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](https://creativecommons.org/licenses/by/4.0/)).

Summary

Metagenomic analyses aim to explore the genomic diversity of communities in specific habitats by processing their DNA sequencing data. This analysis is achieved with specialized bioinformatics tools, which often require previous coding experience. Furthermore, beginners can struggle to build a pipeline from raw data to valuable biological insights. [The Carpentries](#) hosts open lessons used worldwide to analyze specialized datasets for beginners, including [a Data Carpentry curriculum for individuals working with genomics sequencing data](#). However, a lesson addressing the specific challenges associated with metagenomics data and analyses was missing. We created a complete Metagenomics curriculum in [The Carpentries Incubator](#), adapting and expanding on the Data Carpentry genomics curriculum. The curriculum provides an introduction to programming, teaching learners to access and handle metagenomics data, and to run commands with the software needed for completing metagenomics analyses. Content and exercises have been improved based on experience gathered in teaching the curriculum in three 16-hour online workshops. We expect to continue to enhance this lesson, which we hope is helpful as a teaching resource for new instructors in the field, and as a guide for newcomers wishing to perform metagenomic analyses from scratch.

Statement of Need

Bioinformatic tools are now essential to our understanding of biological systems. Open lessons for general-purpose coding languages and specialized topics such as genomics, ecology, and even metagenomics are already available ([Darling et al., 2019](#)), ([Lessons, n.d.](#)), ([Kruchten, 2020](#)). Nevertheless, a complete guide for shotgun metagenomics, assuming no prior knowledge of coding and provided hardware and software solutions, was missing. We introduce a curriculum to fill this gap, which teaches the skills required to conduct a comprehensive metagenomics workflow with Bash and R programming in a pre-installed remote server.

Content and Learning Objectives

The curriculum is comprised of four lessons (Fig 1). It assumes no previous programming experience but expects that students understand basic concepts of molecular biology and microbiology. The first two lessons, introductions to project planning and organization and programming in Bash respectively, are adapted from equivalent lessons in the Data

Carpentry Genomics curriculum ([becker_datacarpentrygenomics-workshop:2019?](#)). The third part includes a brief introduction to R, and the fourth teaches a complete shotgun metagenomics workflow using public data ([Okie et al., 2020](#)).

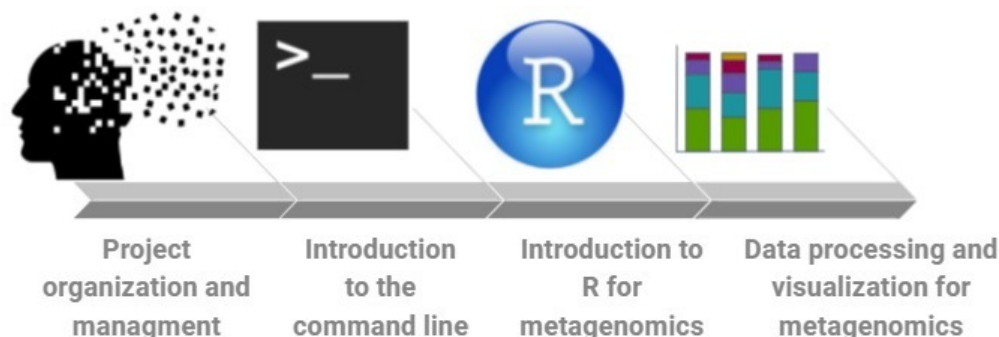


Figure 1: Fig 1.

Fig 1. The four lessons in the curriculum for a metagenomics workshop.

Learning objectives

Project organization and management for metagenomics - Plan, structure, organize, and document metagenomics data, metadata, and bioinformatics workflow. - Access public data on the NCBI sequence read archive (SRA) database.

Introduction to the command-line for metagenomics - Identify the benefits of the command line. - Navigate the file system, create, copy, move, and remove files and directories. - Work with files. - Combine commands and automate repetitive tasks. - Make an organized file structure for a bioinformatics project.

Introduction to R for metagenomics - Familiarize with RStudio and R functionality. - Distinguish the different data types. - Create and manipulate data frames. - Use and find help in R libraries.

Data processing and visualization for metagenomics - Explain the structure and contents of the data and metadata used in these lessons. - Assess the quality of sequencing data. - Trim and filter sequences based on their quality. - Perform a metagenomic assembly. - Obtain Metagenome-Assembled Genomes and check their quality. - Assign and visualize a taxonomy of reads and contigs. - Explore the diversity in a sample and calculate diversity estimates. - Discover more resources for metagenomics projects.

Instructional Design

The workshop includes 16 hours of content with live coding, formative assessment practice, and other supportive elements. During the development of the lessons, we considered three axes of teaching: cloud setup, standardized episodes, and teaching strategies.

1. Cloud setup: We set up remote machines to standardize the learning environment and teaching experience. This setup lowers entry barriers for students without experience in technical installations and provides enough computing power regardless of individual configuration. Students only need a computer with internet access and a terminal program installed. We also enlist the instructions for setting up the remote machines and an alternative installation guide to be used by people who prefer to follow the lessons on a different computer.

2. Standardized episodes: Following The Carpentries instructional design, lessons consist of multiple short *episodes* of content, each prefaced with the questions to be answered and the learning objectives for that section, and ending with concrete, clear messages of the learned content.
3. Teaching strategies
 - Live coding: Programming simultaneously with students gives a practical coding experience with examples of mistakes arising and being solved.
 - Exercises and discussion in small groups: Allows learners to solve problems with peers and incentivizes participation while applying what they learn to strengthen new knowledge and skills. Regular breaks for exercises provide learners and instructors with feedback about progress toward the learning objectives.
 - Content review: At the end and beginning of each session, we ask learners to revisit, list, and explain the content taught in past lessons. We use a collaborative document where simultaneous written and spoken review helps to reaffirm the content.

Teaching Experience

Our teaching team includes undergraduate and graduate students, postdocs, and professors, each with different perspectives that enrich the discussion of the lesson. We found it adequate to have two instructors for the online workshops, plus one helper for every five students. We are pleased to have recruited new helpers from the attendees of previous workshops, enlarging and further enriching our teaching community. To promote an interactive learning experience, we use a collaborative live document and small groups to do exercises and discussions. In this way, there are fewer barriers to student's participation, and they gain practical experience of the different solutions bioinformatic problems can have. We include bonus exercises around the lessons that are not meant to be solved by all the learners but are an opportunity for advanced participants to tackle more challenging tasks. They help to leverage the background and learning speed of the learners and provide opportunities to practice and reinforce the content learned. Another strength of our workshop structure is introducing programming languages followed by a practical example of their usage, so the metagenomic analysis grows while the coding skills learned become meaningful. The lessons are helpful even if the student's primary goal is to learn coding rather than metagenomics. We ask the attendees to fill out surveys before, during, and after the workshops to adapt our teaching strategies to the current and future learning groups. With the results of these surveys, we have been able to improve the curriculum's content, explanations, and exercises.

Story of the Project

The idea to create a lesson comprising all of the steps required to process and analyze metagenomics data arose when one of the development team wanted to learn how to understand metagenomics data but found it overwhelming even to decide where to start. As someone who already used and knew the advantages of the resources offered by The Carpentries, she recognized the absence of a curriculum about metagenomics on this platform. When she learned of The Carpentries Incubator, she reached out to a fellow student who had experience on the subject and started building the episodes. Other students in the same institute, who were also working on metagenomics, started contributing to the pipelines they were learning. Further team members were recruited as helpers after attending a workshop and later began collaborating to develop the episodes.

Acknowledgments

Karina Enríquez Guillén and Rafael Pérez Estrada for testing and improving the lesson content. Diana Oaxaca, Alejandro Pereira Santana, Angélica Ruiz, Brian Bwanya, Ahmed Moustafa and Israel Pichardo, for enriching content and teaching strategy. Angélica Cibrián-Jaramillo, Francisco Barona-Gómez, and Harumi Shimada enabled and promoted the delivery of these lessons among their communities. Developers of Data Carpentry Genomics for their work and ideas. Toby Hodges and Erin Becker for their technical support during the development and delivery of the lessons. Students for their feedback and enthusiasm. We thank UNAM for funding in proyecto PAPIME “Desarrollo de material didáctico de Bioinformática con énfasis en metagenómica para las modalidades presencial y virtual”

Darling, A. E., DeMaere, M. Z., Gaio, D., & Anantanawat, K. (2019). *Data for the gene school: Metagenomics workshop*. Zenodo. <https://doi.org/10.5281/ZENODO.3585993>

Kruchten, A. E. (2020). A curricular bioinformatics approach to teaching undergraduates to analyze metagenomic datasets using R. *Frontiers in Microbiology*, 11, 2135. <https://doi.org/10.3389/fmicb.2020.578600>

Lessons. (n.d.). Retrieved September 12, 2021, from <https://datacarpentry.org/lessons/>

Okie, J. G., Poret-Peterson, A. T., Lee, Z. M., Richter, A., Alcaraz, L. D., Eguiarte, L. E., Siefert, J. L., Souza, V., Dupont, C. L., & Elser, J. J. (2020). Genomic adaptations in information processing underpin trophic strategy in a whole-ecosystem nutrient enrichment experiment. *eLife*, 9, e49816. <https://doi.org/10.7554/eLife.49816>