# Training Biologists in Unix Command-Line Skills: From Competency Framework to a Scalable Self-Paced Implementation

Lucie Khamvongsa-Charbonnier [†] , Robert Aboukhalil [†] , Helene Chiapello [*,†] , Thomas Denecker [†] , Pierre Poulain [†] , Denis Puthier [†] , Olivier Sand [†] , Morgane Thomas-Chollier [†] , Claire Toffano-Nioche [†]

*Article*

# Training Biologists in Unix Command-Line Skills: From Competency Framework to a Scalable Self-Paced Implementation

**Lucie Khamvongsa-Charbonnier [1], Robert Aboukhalil [2], Hélène Chiapello [1,3,\*],
Thomas Denecker [1], Pierre Poulain [4], Denis Puthier [5], Olivier Sand [1],
Morgane Thomas-Chollier [1,6] and Claire Toffano-Nioche [7]**

[1] IFB-core, Institut Français de Bioinformatique (IFB), CNRS, INSERM, INRAE, CEA, 94800 Villejuif, France

[2] OMGenomics Labs, 548 Market St #661003, San Francisco, California 94104, United States

[3] Université Paris-Saclay, INRAE, MaIAGE, 78352 Jouy-en-Josas, France

[4] Université Paris Cité, CNRS, Laboratoire de Biochimie Théorique, 13 rue Pierre et Marie Curie, 75005, Paris, France

[5] Aix-Marseille Université, INSERM, TAGC, MarMaRa Institute, Turing Centre for Living systems, Transcriptomics and Genomics Marseille Luminy (TGML), 13288 Marseille, France

[6] Institut de biologie de l'ENS (IBENS), École normale supérieure, CNRS, INSERM, Université PSL, 75005 Paris, France

[7] Université Paris-Saclay, CEA, CNRS, Institute for Integrative Biology of the Cell (I2BC), 91198 Gif-sur-Yvette, France

**\*** Correspondence: helene.chiapello@inrae.fr

**Abstract**

As the generation of data in the life and health sciences expands rapidly, there is a growing need for professionals and students in these fields to master core bioinformatics skills, particularly those relating to Unix-like systems, most commonly used in bioinformatics. This paper introduces two key contributions to address this need: (1) a detailed Unix competency framework tailored for bioinformatics, which uses the revised Bloom cognitive taxonomy to define progressive and operational skill levels; and (2) a series of interactive online tutorials deployed through Sandbox.bio, an open-source platform for learning bioinformatics that embeds a command-line in the browser, which removes barriers related to software installation and access. Tested in multiple educational and professional settings, the tutorials have received highly positive feedback. This inclusive, sustainable approach provides widespread access to essential bioinformatics skills for students and professionals alike.

**Keywords:** bioinformatics education; e-learning; web assembly; competency framework

## Introduction

As data generation continues to accelerate, the life science and health communities are facing a growing need to master core bioinformatics skills. These communities encompass a diverse range of profiles from university undergraduates to professionals in research and medical institutions, eager to learn new analysis techniques [1–4]. Training such a large and diverse community raises significant challenges, especially regarding content design, trainer availability and access to computational resources [5,6]. Unix, Python and R are commonly identified as core competencies in data science in general and also in many bioinformatics training programs [7,8]. One notable example is the Software Carpentry initiative [9], which has been playing a key role in providing high-quality training and open material to the community (see: https://software-carpentry.org). However, this initiative largely depends on face-to-face training, which creates a bottleneck due to the need for a sufficient number

of trained instructors. It also requires using a commercial cloud (Amazon Web Services) or installing several bioinformatics tools locally (using the command-line interface!).

In recent years, self-paced learning has emerged as a key pillar of education, offering flexible, learner-centered approaches that allow individuals to progress at their own pace, manage their schedules, and revisit foundational concepts as needed [10]. This format is particularly well suited for acquiring prerequisites ahead of in-person courses, managing large student cohorts, and supporting lifelong learning in professional environments. To meet this demand, several platforms—such as DataCamp [https://www.datacamp.com], Katacoda [now closed] and Killercoda [https://killercoda.com]—have contributed to democratizing self-paced learning in fields like data science and programming. However, these solutions remain proprietary, raising concerns about long-term availability (e.g., Katacoda was discontinued in 2022 [11]), and are generally not designed for bioinformatics. More broadly, many free and open educational resources have been developed to support training in computational biology. For instance, the Galaxy Training Network [12,13] is a prominent collaborative platform offering open-source tutorials tailored for both scientists and trainers across a wide range of topics. TeSS, the ELIXIR Training eSupport System, has a dedicated section to discover existing e-learning materials (https://tess.elixir-europe.org/elearning_materials). Additionally, an increasing number of bioinformatics tools now come with dedicated tutorials for educators and beginners, facilitated by the accessibility and collaborative features of platforms like GitHub and GitLab. Yet, they often assume access to sufficient computing infrastructure—whether local machines or remote clusters—and a minimum level of system configuration skills. These technical requirements pose a significant barrier for learners with limited prior experience, especially in institutional settings and countries where access to resources is uneven.

In this context, WebAssembly (Wasm) technologies offer a transformative alternative by enabling programs to run directly within the user's browser. This approach harnesses the user's own computing power, eliminating the reliance on large centralized computational infrastructures. Moreover, software installation and local configuration are handled seamlessly, enabling effortless deployment and a smooth user experience. As most introductory tutorials make use of small datasets for demonstration purposes, a basic web browser serving basic files for a website (e.g., HTML, CSS, Javascript) together with Wasm-compiled programs should be generally sufficient for training. Wasm thus appears as a promising solution to meet the needs of online training for bioinformatics core competencies where complex toolchains and system dependencies often hinder early learning. For instance, JupyterLite [https://jupyterlite.readthedocs.io/en/latest/], a Wasm implementation of JupyterLab, allows running a fully-fledged JupyterLab in the browser without the need of any external server.

Relying on this innovative technology, Sandbox.bio [https://sandbox.bio] represents a forward-looking step in bioinformatics education. At the initiative of one of the authors of this article (RA), it has now become an open-source online platform specifically designed for bioinformatics training. This platform leverages the use of the Wasm technology to provide self-paced guided tutorials for popular bioinformatics programs such as BLAST, SeqKit or fastp. Sandbox.bio also provides virtual playgrounds where learners can experiment with Unix tools such as Grep, Awk or Sed.

Despite the consensus that Unix is a fundamental skill in bioinformatics, the specific competencies required in the field of bioinformatics are rarely defined beyond a generic level like "basic", "intermediate" or "advanced", rarely attached to any precise operational competencies. For example, does a biology student learning Unix need to know all or only some of the Unix commands related to data manipulation (like, grep, cut, count, awk...)? Does he/she need to learn to write a complete Bash script or rather only need to understand how a given script is working and be able to adapt it to his/her needs? It has become crucial to go beyond generic levels to truly unlock self-paced learning, and globally facilitate learning of Unix competencies for the life sciences.

In this work, we address two issues: Firstly, we define Unix competency levels in bioinformatics that go beyond the three generic basic, intermediate and advanced levels. Secondly, we propose sustainable online tutorials to help biologists master these competencies. We propose two

complementary resources. The first is a core Unix competency skill framework that details and organizes each individual Unix competency in coherent and progressive categories. This framework can be used for various purposes, such as self-assessing student competencies, designing a training curriculum based on learner profile, or specifying target competencies to be acquired upon completion of the training. Secondly, we propose a series of interactive online tutorials based on the WebAssembly technology and hosted on the open-source sandbox.bio platform. This enables access without installation and supports a large number of trainees.

These two resources can be used independently, yet they were designed to be complementary, with the tutorials directly covering the competencies outlined in the skill framework. Designed for biologists, these resources enable the progressive learning of bioinformatics with biology-specific examples. Hosted by the sandbox.bio platform, they are inexpensive in terms of computing resources and can be deployed on a large scale without posing any beginner-related risk to computing infrastructures, such as breaking something on shared servers.

## Material and Method

*Unix Skill Levels in Bioinformatics*

To define the skill levels, we first described Unix relevant skills for bioinformatics within the framework of the revised Bloom cognitive taxonomy framework [14]. Bloom's taxonomy classifies educational objectives into six different categories: Remember, Understand, Apply, Analyze, Evaluate, and Create. We used these categories to compile a list of increasing complexity abilities within the competency framework. These abilities were also grouped into themes, ranging from understanding what a command-line interface is, to running a genomics analysis software in the terminal.

*Tutorial Implementation*

Based on the Unix skill levels, we created step-by-step tutorials to guide learners in acquiring competencies defined for each skill level. Each tutorial consists of a couple of steps interspersed with quizzes, and covers at least one of the categories from Bloom's taxonomy. The tutorials were written in Markdown format and parametrized with a JSON file.

*Tutorial Deployment*

The tutorials were deployed on our own instance of the Sandbox.bio platform [https://github.com/sandbox-bio/sandbox.bio]. This platform offers command-line emulation based on the v86 project [https://github.com/copy/v86], a Wasm-compiled x86 emulator that runs in the browser. We built the Docker image that powers the command-line interface on an Ubuntu 22.04.5 LTS virtual machine with 8 GB of RAM and 2 CPUs, and Node.js was used as the web server. We adapted the code so that the folder tree mimics the one of the ELIXIR-FR/IFB high-performance computing infrastructure. This ensures that users learning on these tutorials are well-prepared for moving afterwards on the infrastructure.

Once deployed, each tutorial is accessible as a "*scenario*" on a public web page hosted by the French Institute of Bioinformatics (ELIXIR-FR/IFB), which is part of the European Elixir infrastructure.
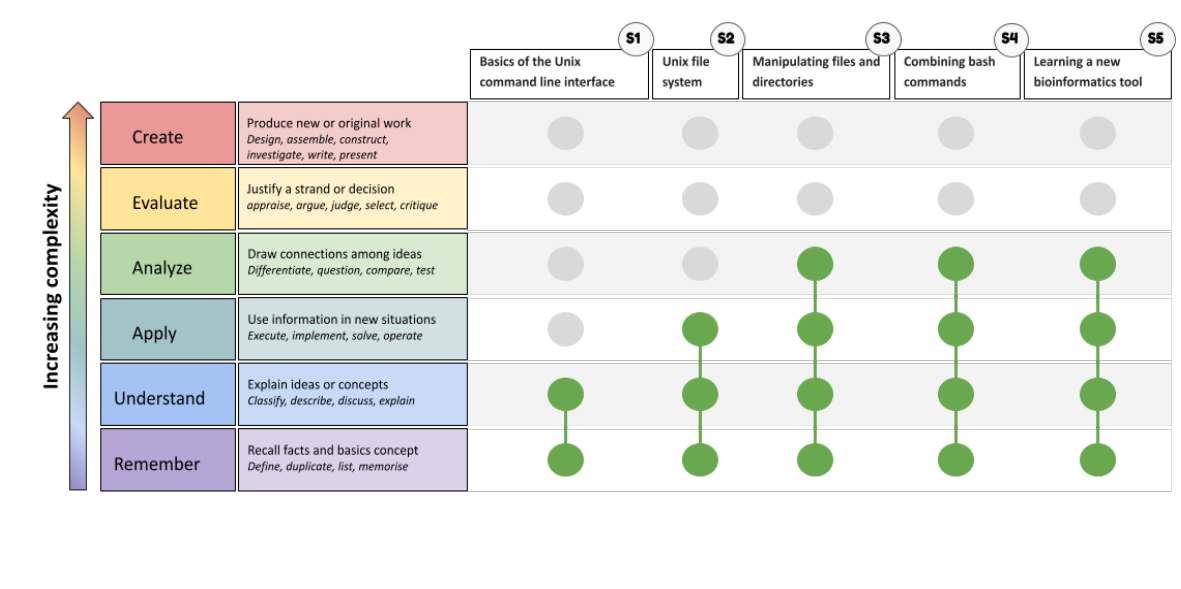
## Results and Discussion

*A Unix Competency Matrix for Learning Bioinformatics*

We designed a competency matrix that thematically groups bioinformatics Unix competencies and uses Bloom's taxonomy to classify the required skills (see Figure 1). The level of difficulty increases, with the skills becoming cumulative: those acquired at one level are necessary to reach the next. We have defined four level of competencies (see Figure 2): U1 (terminal, shell and Unix

commands), U2 (structure and operation of a command in the file system), U3 (advanced exploration of the file system, and workspace), and U4 (motifs, redirection, pipe, flux). Figure 2 includes examples of the learning outcomes for each level. The full competency matrix is freely available in Zenodo in French [15] and in English [16].



**Figure 1.** Application of Bloom's revised taxonomy activities to define groups of Unix competencies of increasing difficulty. The left-hand side of the figure shows the revised Bloom's taxonomy. The columns at the top describe the groups of Unix skills that are mandatory for bioinformatics, classified by theme and ordered by increasing difficulty. These groups are named S1 to S5, in reference to the scenarios developed in the sandbox.bio tutorials.



**Figure 2. The four competency levels and learning outcomes for bioinformaticians (U1 to U4)**. The table presents examples of learning outcomes for each level; refer to the complete documents to view the full extent [15,16]. U5 to U7 refers to three additional levels that will be defined in future versions.

*Online Tutorials for Practical, Self-Paced and Autonomous Learning of Unix*

We used the Unix competency matrix to design five Unix tutorials. Table 1 illustrates the complementary nature of the competency framework and the tutorials, which are deeply

interconnected yet can be used independently for specific purposes. The competency framework structures skill development, while the interactive tutorials provide accessible, hands-on practice. Together, they form a coherent learning ecosystem. Table 1 summarizes the key features of this complementarity, regarding format, target audience, use cases, and access.
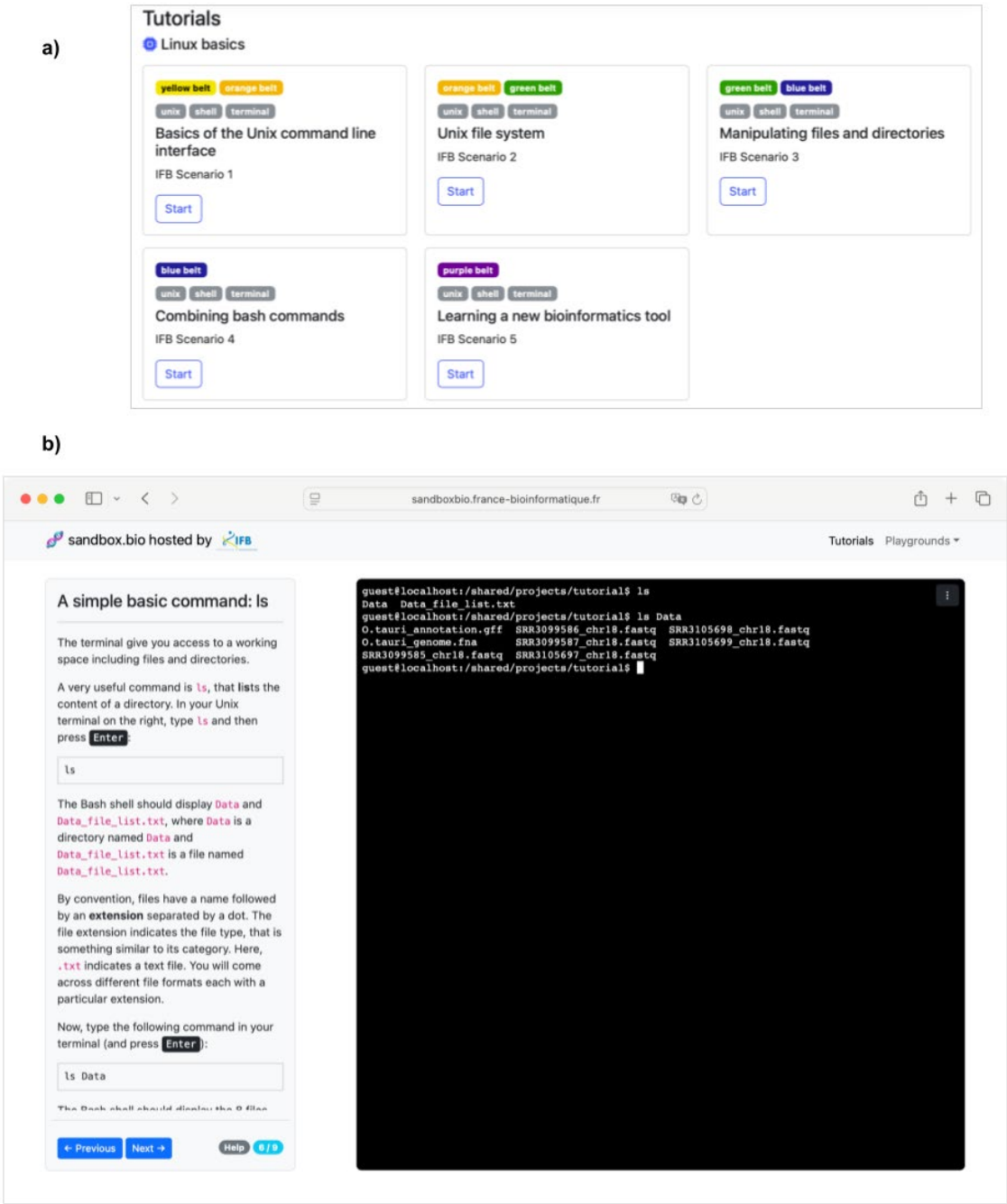
**Table 1. Relations between the Unix competency framework and the interactive online tutorials.** The table illustrates the complementarity of the two resources regarding description, format, target audience, use cases, and access.

|  | Unix competency framework | Interactive online tutorials |
|---|---|---|
| **Description** | A themed, judo-belt-inspired tool to guide the progressive learning of core Unix skills in bioinformatics | Open-access Unix tutorials focused on common bioinformatics applications, flexible and accessible without installation |
| **Format** | Textual framework / PDF | Web-based interactive |
| **Target audience** | Trainers, instructional designers and learners | Learners (beginners to intermediate) |
| **Use cases** | Curriculum design, learner self-assessment, certification | Self-training, classroom support, flipped classroom formats |
| **Access** | Downloadable / Shareable via Zenodo | Online (sandbox.bio) |

Unix tutorials are freely available on github [https://github.com/IFB-ElixirFr/sandboxbioscenarios/tree/master] and via a dedicated Sandbox.bio instance at the French Institute of Bioinformatics (ELIXIR-FR/IFB): https://sandboxbio.france-bioinformatique.fr. Figure 3 shows an example of the website interface and tutorial content. Each tutorial includes pedagogical support, with explanations of concepts and commands on the left, and a Linux terminal for running Unix commands in a secure environment on the right side.

The implemented scenarios have been tested in multiple real contexts with the following settings:

- Pre-class activities in autonomy for professional bioinformatics training: "Introduction to the processing of genomic data obtained by high-throughput sequencing" at the French Bioinformatics Institute (EB3I 2023 and EB3I 2024, with 39 trainees in each session), a one-week intensive training session for life scientists.
- Pre-class activities in autonomy for a professional diploma at Université Paris Cité: "Production, analysis, and valorization of biological omics data" (DU omiques 2023 and DU omiques 2025, with 14 trainees in each session).
- Pre-class activities in autonomy for biologists and interns (14 participants) at the Institute for Integrative Biology of the Cell (I2BC) at Université Paris-Saclay
- Self-assessment of Linux competencies by 30 students of the IMALIS life science master's program at the École Normale Supérieure (Paris).
- Self-assessment of Linux competencies by 40 students of the Polytech graduate school of engineering (specialty in Biological Engineering) at Aix-Marseille University.

**Figure 3. Unix tutorials released at the French Institute of Bioinformatique (IFB/ ELIXIR-FR)** (a) List of the five tutorials currently available on the IFB sandbox.bio instance (b) An example of the user interface when running the first tutorial "Basics of the Unix command line interface" (also named Scenario 1).

We evaluated the online tutorials through a survey with a group of EB3I trainees after the 2023 EB3I training: 90% of participants found the tutorials clear, and 95% considered them to be of an appropriate length. In the EB3I 2024 edition, these figures increased to 97% and 92%, respectively, indicating overall positive feedback of the tutorials. Participants quoted: "I loved it, very instructive.", "Very practical, it allows me to learn about command lines and Bash.", " [This online training] is essential because it lays the necessary foundations for the week. I found this pre-training very appropriate and well-guided, with the option of typing command lines in parallel." One key valuable advantage of our training resources is that they combine a text tutorial with a terminal that enables trainees to practice in a safe environment and at their own pace.

In light of the highly positive responses to this evaluation of tutorials, which can be used by students and professionals alike at their convenience without the need to install any software, we plan to extend their use to other scenarios. For example, we are considering a tutorial to help users access the ELIXIR-FR/IFB high-performance computing cluster.

Interestingly, our work meets several of the grand challenges in bioinformatics education and training recently pointed out by Işık *et al.* [5]: GC1: Identifying fundamental knowledge and skills, GC2: supporting lifelong training, GC6: Practicing inclusivity and equity in bioinformatics education and GC7: Ensuring the sustainability of bioinformatics education.

At this stage, the developed training resources are limited to the basic Unix core skills required for using bioinformatics tools. While this alone does not provide a full bioinformatics training, it serves as an essential first step toward becoming a bioinformatics user, as defined in [6]. This approach meets a significant demand within the Life Sciences, Agronomy and Medical Research communities. Moreover, additional tutorials and resources are available for those wishing to further develop their Unix skills. For instance, the sandbox.bio platform and other facilities offer additional tutorials, such as "*How to write a bash script*" and "*Data exploration with awk*" (available at https://sandbox.bio/tutorials). At I2BC-University Paris-Saclay, GameShell (https://github.com/phyver/GameShell) was used to go deeper after the Sandbox tutorials. Although this tool requires local installation, it guides participants through Unix commands with a dungeon discovery adventure. Going forward, we intend to continue developing free online tutorials to promote the large-scale development of bioinformatics skills.

## References

1.      Attwood TK, Blackford S, Brazas MD, Davies A, Schneider MV. A global perspective on evolving bioinformatics and data science training needs. Brief Bioinform [Internet]. 2017 Aug 29 [cited 2025 Aug 20];20(2):398–404. Available from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6433731/

2.      Wilson Sayres MA, Hauser C, Sierk M, Robic S, Rosenwald AG, Smith TM, et al. Bioinformatics core competencies for undergraduate life sciences education. PLoS ONE [Internet]. 2018 June 5 [cited 2025 Aug 20];13(6):e0196878. Available from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5988330/

3.      Mulder N, Schwartz R, Brazas MD, Brooksbank C, Gaeta B, Morgan SL, et al. The development and application of bioinformatics core competencies to improve bioinformatics training and education. PLoS Comput Biol [Internet]. 2018 Feb 1 [cited 2025 Aug 20];14(2):e1005772. Available from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5794068/

4.      Perkel JM. Five reasons why researchers should learn to love the command line. Nature [Internet]. 2021 Feb 4 [cited 2025 Sept 9];590(7844):173–4. Available from: https://www.nature.com/articles/d41586-021-00263-0

5.      Işık EB, Brazas MD, Schwartz R, Gaeta B, Palagi PM, Van Gelder CWG, et al. Grand challenges in bioinformatics education and training. Nat Biotechnol [Internet]. 2023 Aug [cited 2025 Aug 20];41(8):1171–4. Available from: https://www.nature.com/articles/s41587-023-01891-9

6.      Schneider MV, Watson J, Attwood T, Rother K, Budd A, McDowall J, et al. Bioinformatics training: a review of challenges, actions and support requirements. Brief Bioinform [Internet]. 2010 Nov 1 [cited 2025 Aug 20];11(6):544–51. Available from: https://academic.oup.com/bib/article-lookup/doi/10.1093/bib/bbq021

7.      Welch L, Brooksbank C, Schwartz R, Morgan SL, Gaeta B, Kilpatrick AM, et al. Applying, Evaluating and Refining Bioinformatics Core Competencies (An Update from the Curriculum Task Force of ISCB's Education Committee). PLOS Comput Biol [Internet]. 2016 May 13 [cited 2025 Aug 20];12(5):e1004943. Available from: https://dx.plos.org/10.1371/journal.pcbi.1004943

8.      Welch L, Lewitter F, Schwartz R, Brooksbank C, Radivojac P, Gaeta B, et al. Bioinformatics Curriculum Guidelines: Toward a Definition of Core Competencies. PLoS Comput Biol [Internet]. 2014 Mar 6 [cited 2025 Aug 20];10(3):e1003496. Available from: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3945096/

9.      Wilson G. Software Carpentry: Getting Scientists to Write Better Code by Making Them More Productive. Comput Sci Eng [Internet]. 2006 Nov [cited 2025 Aug 20];8(6):66–9. Available from: http://ieeexplore.ieee.org/document/1717319/

10. Lui RWC, Dik MNY, Lee PTY. Evaluation of Online Learning Platforms for Interactive Self-paced Learning of Container-based Software Development. In: 2022 IEEE International Conference on Teaching, Assessment and Learning for Engineering (TALE) [Internet]. Hung Hom, Hong Kong: IEEE; 2022 [cited 2025 Aug 20]. p. 54–8. Available from: https://ieeexplore.ieee.org/document/10148308/

11. Baldwin L. Leveraging Katacoda technology exclusively within O'Reilly And the decision to shutter katacoda.com [Internet]. 2022. Available from: https://www.oreilly.com/online-learning/leveraging-katacoda-technology.html

12. Batut B, Hiltemann S, Bagnacani A, Baker D, Bhardwaj V, Blank C, et al. Community-Driven Data Analysis Training for Biology. Cell Syst [Internet]. 2018 June [cited 2025 Aug 20];6(6):752-758.e1. Available from: https://linkinghub.elsevier.com/retrieve/pii/S2405471218302308

13. Hiltemann S, Rasche H, Gladman S, Hotz HR, Larivière D, Blankenberg D, et al. Galaxy Training: A powerful framework for teaching! Ouellette F, editor. PLOS Comput Biol [Internet]. 2023 Jan 9 [cited 2025 Aug 20];19(1):e1010752. Available from: https://dx.plos.org/10.1371/journal.pcbi.1010752

14. Anderson LW, Krathwohl DR. A taxonomy for learning, teaching, and assessing: A revision of Bloom's taxonomy of educational objectives. 2001. (New York: Longman).

15. Chiapello H, Denecker T, Khamvongsa Charbonnier L, Poulain P, Puthier D, Sand O, et al. Matrice de compétence Unix pour l'apprentissage de la Bioinformatique [Internet]. Zenodo; 2024 [cited 2025 Aug 21]. Available from: https://zenodo.org/doi/10.5281/zenodo.12104694

16. Chiapello H, Denecker T, Khamvongsa Charbonnier L, Poulain P, Puthier D, Sand O, et al. Unix Competency Framework for Learning Bioinformatics [Internet]. Zenodo; 2025 [cited 2025 Sept 9]. Available from: https://zenodo.org/doi/10.5281/zenodo.17084249