

# BIOL 432 – Syllabus

## Introduction to Computation and Big Data in Biology

### Course Information

Winter Semester, 2023. January 9 to April 3

3.0 Credits

Hybrid (In-Person and Online)

**Pre-requisites:** BIOL 343

**Lectures** Mondays, 2:30-4:30, Ellis 319

**Tutorials** Wednesdays, 8:30-9:30, Ellis 319

### Instructor

*Name:* Dr. Robert I. Colautti (he/him)

*Office Address:* 4325 BioSciences, 116 Barrie St.

*Office Hours:* Mondays, 4:30-5:30

*Phone:* 613-533-2353; Please use only during office hours.

*Email:* [robert.colautti@queensu.ca](mailto:robert.colautti@queensu.ca)

*About Me:* My teaching takes a student-centered approach that is supportive of diverse learners. I set high expectations and provide extensive resources so that student learning extends beyond the timeline course. I try to emphasize a growth mindset that is focused on effort, personal development, and quality of work, rather than mastery or excellence. I try to teach students how to recover and learn from failure, which I believe is essential for a successful career. My teaching philosophy draws on two decades of research and mentorship experience, resulting in primary research in top journals (e.g. Science, PNAS, PRSB) and dozens of former students employed in the public and private sectors or continuing to advanced degrees. I have a broader perspective than many biologists, drawing from an MSc in aquatic ecology from the University of Windsor and a PhD in quantitative genetics at the University of Toronto, followed by bioinformatics and computational biology research at Duke University (North Carolina), the University of British Columbia, and the University of Tuebingen (Germany). I am involved in a wide range of research projects, from COVID-19 and Lyme disease to plant ecology and rapid evolution. As a faculty member in the Biology Department at Queen's University, my goal has been to develop courses that cover everything I wish I had known when I was starting out. The result of this effort is the course material for BIOL 343, BIOL 432, BIOL 860 and BIOL 812.

*Website:* <https://ecoevogeno.org/>

*Twitter:* @ColauttiLab

### Important University Dates

Key dates (first day of class, tuition due date, last day to add/drop courses) are important to your academic success. Please find them at [Important Dates](#).

### Welcome to BIOL 432!

This course is designed for students who want to pursue a career in bioinformatics and/or computational biology, building on the foundations of BIOL 243 and BIOL 343. This is taught through custom tutorials written specifically for this course.

We begin with a 'crash course' that reviews fundamentals of the R programming language taught in previous courses. We then explore the Python coding language, emphasizing key similarities and

differences with R that will make it easier for you to understand how to program across environments. The 'crash course' also includes an introduction to high performance computing with Linux and bash programming, after which you will gain access to high-performance servers at Queen's Centre for Advanced Computing (CAC).

Once these basic skills are established, we focus on specific applications in DNA/RNA analysis, phylogenetics, metagenomics/metabarcoding, genome assembly, and transcriptomics. Regular assignments and quizzes will help you develop your coding skills. In addition to individual assignments, you will work in groups for much of this course, ending with a final project in which your group will address an important biological question of your choosing.

The best way to learn coding is with extensive practice, trial, and failure. For this reason, there are no major exams. Instead, you are assessed through regular assignments and quizzes designed to encourage skill development and problem-solving. In addition, you are **STRONGLY ENCOURAGED** to find opportunities to practice coding wherever possible. It will take longer at first, but it will save a lot of time in the long run. Finally, be prepared to get frustrated – you will make many mistakes and most of your coding time will be debugging and searching for answers on the internet. It is important to know that **this is COMPLETELY NORMAL** and self-directed research to solve coding problems is perhaps the most crucial skill you will learn in this course.

## Quotes from previous students:

*I thought I would be underprepared but it was well-tailored to students with limited coding/computer science experience*

*I have no clue why the assignments were only open for around 48hours. As a student this was a pain and definitely lead to me rushing through some of them*

*I worked very hard in this course, and I'm pleased to report that it paid off--I feel fairly confident in applying the skills I learned here. I hadn't been looking forward to it when I signed up (I just needed a 400lvl to graduate) but I ended up enjoying it a lot. By the end, I was actually starting to agree with Maria that "this is actually a really cool assignment". I especially liked the Dragon Phylogeny (I sent that report to my dad, he liked it too).*

*Overall this course was really helpful! I feel that I have learned a lot. The self-tutorials were very useful because they let me work through the code at my own pace, and I appreciate the weekly quizzes to help solidify what I've learned. This course was so much more fun and interesting than I had originally thought (since coding is not my strong suit) and I feel like I have accomplished a lot! Thanks for teaching it!!*

## Equity, Diversity and Inclusivity Statement

We recognize equity and diversity are central to our educational mission and standards of excellence. We are working to dismantle direct, indirect, and systemic discrimination that still exists within our institutional structures, policies and practices -- and in our community. These take many forms and work to differentially advantage and disadvantage persons across social identities such as race, ethnicity, disability, gender identity, sexual orientation, faith and socioeconomic status, among other examples. As students and educators, we have important roles to play to identify and address systemic discrimination for the benefit of all.

## Land Acknowledgement

Queen's is situated on traditional Anishinaabe and Haudenosaunee territory. We are so privileged and grateful to be able to live, learn and play on these lands. Therefore, let us be mindful of the what we learn while on these lands, and how we might apply our newfound skills and knowledge for the benefit of all.

## Expectations

### For Instructors & Teaching Assistants

As course instructor, I am responsible for developing and editing the course material, which was written specifically for this course and for you, the student interested in bioinformatics and computational biology. I will ensure that all relevant course material is available online and released on a weekly basis so that you know where you should focus your limited time and attention. The course content is a work in progress, so I welcome any feedback on this material, from small spelling/grammar errors to general suggestions for improvement.

To accommodate differences in learning styles, I will make the main content available in complementary forms including an online textbook and pre-recorded videos with annotated scripts. The textbook will lean heavily on a tutorial style, with step-by-step instructions that are reiterated in the online videos. The videos and textbook are designed to be complementary, with overlap emphasizing important skills and techniques. I will also be responsible for recording and posting in-person lectures for those who are not able to attend, or wish to review, the synchronous sessions.

I will make mistakes in these tutorials. Everyone makes mistakes, and coding is particularly prone to error, especially when there are distractions. I will use these opportunities to demonstrate how to troubleshoot errors by carefully reading the warning messages and running smaller subsets of code to identify where the problems lie. **Learning how to troubleshoot mistakes is one of the most important skills you can learn in this course.**

The entire teaching team (instructor + teaching assistants) is committed to maintaining a healthy and inclusive learning environment. Just as we make mistakes, we recognize that mistakes and errors are an important part of your learning process too. We respect and value students who are not afraid to take risks or try things that might be 'wrong'. Above all, we value students who are not afraid to fail. We will use frequent assignments and testing to limit the impact of mistakes on your final grade. We will provide timely feedback – usually within two weeks – so that you can improve your grade on future assignments.

We will communicate twice per week during lecture and tutorial. Lectures will cover only part of the scheduled time allotment, so that there will be ample time available for questions or assistance.

### For Students

It is expected that you will attend weekly lectures and tutorials, though we understand this may not be possible for everyone, all the time, particularly in the post-COVID era. Therefore, everything you need to succeed in this course will also be available online.

You are expected bring a laptop capable running Windows, MacOS, or Linux programs, and you must be able to access Queen's wireless network during lecture and tutorial sessions. When working in class or following lectures at home, **it is very important that you code along in real time.** The only way to effectively learn to program is to practice, and you are expected to practice as much as possible!

You are expected to complete quizzes online, which are administered at the beginning of lecture or tutorial. Assignments are also submitted online and generally due with 48-72 hours of being posted online. You will be prepared for these short deadlines, which are essential to reinforce what you learn in the lectures and online self-tutorials. You must complete quizzes and assignments alone, without communicating with other students. **Any attempts to communicate about quiz or assignment answers will be treated as a breach of academic integrity.** Plan to devote 3-5 hours to learning the lecture material and 10 hours to complete the assignments.

You are expected to check the course website regularly (or use alerts) to keep track of deadlines. **Late assignments are scored as zero** (but see below regarding accommodations).

Any questions or concerns about the course should be raised in lecture or tutorial, or privately during weekly office hours (no appointment needed). Email is generally not an effective tool for course material, and questions that can be addressed in person will not receive an email response. However, email is encouraged for urgent issues (e.g., broken website links or other time-sensitive errors).

### For Interactions

You will have regular interactions with the teaching team (TAs, instructor) and with your classmates. In all interactions, you are expected to be respectful and always behave with integrity, both in face-to-face interactions and when engaging online.

This course will also involve group-based activities that will require communication outside of the classroom. You are responsible for maintaining contact and collaborating with all members of your group in a respectful and timely manner. Remember that other members of your group may not have the same resources or privileges and may need some flexibility or accommodation. **Developing skills to collaborate effectively within a diverse group of peers is an important learning outcome of this course.**

## Course Learning Outcomes

Students completing this course shall be able to:

1. Design and implement a project management strategy that is **OPEN** and **REPRODUCIBLE**.
2. Write custom scripts to curate, merge, subset, reformat, and parse large biological datasets
3. Analyze and interpret 'big data' formats in biology (e.g. CSV, FASTA, FASTQ, SAM, BED, BAM, KML, XML, BMP, PNG, SVG, SHP) to address biological hypotheses.
4. Write clean and coherent code that combines R, Python, and Linux/bash techniques, including analysis pipelines run on remote servers maintained by Queen's Centre for Advanced Computing (CAC).
5. Use regular expressions to modify biological data files (e.g., automated error correction, file conversion, and data extraction)
6. Use Git with GitHub to collaborate with peers on large coding projects.

## Course Materials

The following texts are required for this course:

*R Crash Course for Biologists* by Robert I. Colautti. This book is available for free via the course website and reviews key aspects of R. This book adopts a self-tutorial style, and it is important that you take the time to physically type out the commands in your computer. The simple act of typing is critical to develop coding skills. Sometimes you will not get the same output, and that's a good opportunity to learn how to troubleshoot typos and other errors. You should already be familiar with most of this content, so you should be able to work through it quickly.

*Bioinformatics for Biologists* by Robert I. Colautti. This book is also available for free via the course website and forms the basis for bioinformatics lectures. As with the R Crash Course, this book adopts a self-tutorial style and you should physically code along on your computer.

In addition to the above required reading, I highly recommend the following books if you are serious about pursuing a career in bioinformatics and/or computational biology:

*Practical Computing for Biologists* by Haddock & Dunn (Oxford University Press). This book also follows a self-tutorial style but leans more heavily on Python

*Bioinformatics Data Skills* by Vince Buffalo (O'Reilly). This book also follows a tutorial style with more focus on examples from genomics with high-performance computing. It includes a combination of R, Python and Linux.

## Course Timeline

The following is the planned timeline for the course, however an updated version is available on the course website. Note that there is a **quiz each week**, which is administered at the beginning of class or tutorial. There is also an **assignment due each week**, which is typically posted at the end of lecture and due within 48-72 hours.

	Topic
<b>Week 01</b>	R-eview Pt 1: Base R, visualizations, and data wrangling
<b>Week 02</b>	R-eview Pt 2: Statistical models in R
<b>Week 03</b>	R-egex (regular expression)
<b>Week 04</b>	Collaborating with Git and GitHub in R Studio
<b>Week 05</b>	Introduction to Python
<b>Week 06</b>	Introduction to Linux & SLURM for high-performance computing
<b>Week 07</b>	DNA sequence data and alignments
<b>Week 08</b>	The Dragon Phylogeny
<b>Week 09</b>	Metabarcoding
<b>Week 10</b>	Genome assembly
<b>Week 11</b>	Transcriptomics
<b>Week 12</b>	Final GROUP presentations

## Suggested Time Commitment

Each week, you should spend 1 to 3 hours reviewing lecture videos, 3 to 6 hours working through the tutorials, and 3 to 12 hours completing assignments. In the last 5 weeks of class, you should devote an additional 10 hours per week to your final group assignment. Note the wide range of expected time to complete the assignments, which depends on an unpredictable amount of time needed for troubleshooting. It is strongly recommended that you budget 12 hours or more for each assignment, to ensure that you do not miss the submission deadline.

## Assessment

40%	Weekly Assignments
10%	Weekly Quizzes
10%	Final Project – Proposal (group)
10%	Final Project – Report (group)
10%	Final Project – Code (group)
20%	Participation and peer evaluation

## Essential Requirements and Flexibility to Succeed

The workload in this course is demanding because frequent and regular practice is essential to develop competence as a data scientist. However, we also understand that it may not be possible for every student to commit 20+ hours every week to this course. Therefore, we have incorporated a “Universal Design for Learning” that adds flexibility for students who require accommodations, without the need to register with QSAS or request permission from the teaching team.

Instead of traditional accommodations, we account for unforeseen circumstances in our course design. To do this, we will remove the two lowest weekly quiz scores and the two lowest assignment scores before calculating the final grade. This means that students can miss up to two weekly quizzes and two weekly assignments without penalty. Additionally, students may take up to 1 additional day to complete a quiz and 3 additional days to complete an assignment.

## Grading Scheme and Grading Method

All components of this course will receive numerical marks, weighted by the percentage shown in the “Assessment” section, above. The final grade you receive for the course will be derived by converting your numerical course average to a letter grade according to Queen’s Official Grade Conversion Scale:

### **Queen’s Official Grade Conversion Scale**

<b>Grade</b>	<b>Numerical Course Average (Range)</b>
A+	90-100
A	85-89
A-	80-84
B+	77-79
B	73-76
B-	70-72
C+	67-69
C	63-66
C-	60-62
D+	57-59
D	53-56
D-	50-52
F	49 and below

## Questions About the Course and Contacting the Teaching Team

Coding involves a lot of trial-and-error that can be frustrating for students new to the discipline. First, know that this is completely normal, even to seasoned data scientists. A very common and effective approach to solving errors or other problems is to search Google or Stack Overflow. Often, simply copying and pasting an error into an online search will produce a helpful link. Very often you can just type ‘How do I X in R’ (or the R package name like ggplot2, dplyr) into Google and look for links to similar questions answered on the Stack Overflow website. In addition, we will generally leave ample time at the end of lectures and tutorials. You are strongly encouraged to ask the question in lecture or tutorial so that all students can benefit from the answer. Any private questions or issues can be discussed during office hours (no appointment necessary). Email is not an effective mode of communication in this course, except to notify the teaching time of potential errors in the online material.

## Course Announcements

News and general announcements are posted on the home page of the course website, and new content is released to the course website every week (e.g., lecture topics, assignments). I strongly encourage you to check the website after each lecture.

## Course Feedback

Course feedback is solicited throughout the semester. Your feedback is very important to the teaching team and may be used to adjust the course both in the future and the present. In addition to the feedback that we ask for during lectures and tutorials, we welcome any suggestions for improvement that you would be willing to provide (email is a good format for this).

## Academic Accommodations and Considerations

As noted above, this course uses a Universal Design for Learning philosophy, adding flexibility in deadlines and grading for students who need accommodations. Any additional accommodations should be handled through Queen's Student Accessibility Services (QSAS):

<https://www.queensu.ca/studentwellness/accessibility-services>

## Academic Integrity

Queen's students, faculty, administrators and staff all have responsibilities for upholding the [fundamental values of academic integrity](#); honesty, trust, fairness, respect, responsibility and courage. These values are central to the building, nurturing, and sustaining of an academic community in which all members of the community will thrive. Adherence to the values expressed through academic integrity forms a foundation for the "freedom of inquiry and exchange of ideas" essential to the intellectual life of the University (see the [Senate Report on Principles and Priorities](#)).

Students are responsible for familiarizing themselves with the regulations concerning academic integrity and for ensuring that their assignments and their behaviour conform to the principles of academic integrity. Information on academic integrity is available in the Arts and Science Calendar (see [Academic Regulation 1](#)), on the [Arts and Science website](#), and from the instructor of this course. Departures from academic integrity include plagiarism, use of unauthorized materials, facilitation, forgery, use of forged materials, contract cheating, unauthorized use of intellectual property, unauthorized collaboration, failure to abide by academic rules, departure from the core values of academic integrity, and falsification, and are antithetical to the development of an academic community at Queen's. Given the seriousness of these matters, actions which contravene the regulation on academic integrity carry sanctions appropriate to the severity of the departure that can range from a warning or the loss of grades on an assignment to the failure of a course to a requirement to withdraw from the university.

**Plagiarism** is a form of cheating and includes copying code written by students. There is no 'right answer' for the assignments in this class – there are often many potential coding solutions. You will also develop your own coding style, which will make it obvious when code has been copied. To avoid potential for plagiarism, ALWAYS COMPLETE ASSIGNMENTS ON YOUR OWN. As a bonus, you will learn to code better. On the other hand, it is completely fine to ask others to help you troubleshoot an error message or help you figure out why your code isn't working properly. If you become aware of anyone trying to share or solicit code for the assignments, please point them to this passage and inform the teaching team immediately.

## Copyright of Course Materials

Course materials created by the course instructor, including the textbooks, online tutorials, slides, presentations, quizzes, assignments, and other similar course materials, are the instructor's intellectual property. It is a departure from academic integrity to distribute, publicly post, sell or otherwise disseminate an instructor's course materials or to provide an instructor's course materials to anyone else for distribution (including note sharing sites), posting, sale or other means of dissemination without the instructor's express consent. A student who engages in such conduct may be subject to penalty for a

departure from academic integrity and may also face adverse legal consequences for infringement of intellectual property rights.

## Technology Requirement

You must have a laptop computer with internet access to participate in this course. Before attending the first lecture, you should install the latest versions of the following software:

- The R programming environment (free): <https://www.r-project.org/>
- R Studio Desktop (Open Source Edition, free): <https://www.rstudio.com/products/rstudio/#rstudio-desktop>
- Open R Studio and run the following lines in the terminal and press enter after each. NOTE: this will install some of the R packages that we use in the course. It may take several minutes to install each one. Be sure to type each line EXACTLY:
  - a. `install.packages("ggplot2")`
  - b. `install.packages("tidyverse")`
- Install Python (free): <https://www.python.org/downloads/>
- If you are using a Windows laptop, Install MobaXTerm (MacOS and Linux users do not need this): <https://mobaxterm.mobatek.net/>
- Install Git (free) <https://git-scm.com/book/en/v2/Getting-Started-Installing-Git>
- Sign up for a (free) GitHub Account: <https://github.com/>