

Cultivation of Enhanced Bioinformatic-Specific Pedagogical Manipulatives, Interventions, and Professional Development

by

Marcus D. Sherman

A dissertation submitted in partial fulfillment
of the requirements for the degree of
Doctor of Philosophy
(Bioinformatics)
in the University of Michigan
2022

Doctoral Committee:

Associate Professor Ryan E. Mills, Chair
Associate Professor Alan Boyle
Associate Professor Rajesh Mangrulkar
Professor Maureen Sartor
Professor Patrick Schloss
Assistant Professor Joshua Welch

Marcus D. Sherman

mdsherm@umich.edu

ORCID iD: 0000-0002-0243-4609

© Marcus D. Sherman 2022

Dedication

God. Wife. Everything else.

Acknowledgements

“Do nothing from selfish ambition or conceit, but in humility count others more significant than yourselves. Let each of you look not only to his own interests, but also to the interests of others” (Philippians 2:3-4). The ancient Greek word *propempo* means to set one forward and fit them out with the requisites for a journey: in other words, to set someone up for success. This section is not about me. It is about those around me. It is about those who shaped me and shaped this. This section is about *propempo*.

First and foremost, I *need* to acknowledge my Lord and Savior, Jesus Christ—the Son of God. I never could have fathomed the possibility of the present when I surrendered all to Him not so long ago. God has given me strength and patience when I had none. He gave me grace and mercy when I deserved none. He even made the paths I perceived as broken and endless to be straight and good. Without Him, none of this was possible. Without Him, none of this is possible.

God also made me better by pairing me with my wife, Julie. Through her I was able to do everything from celebrate my first “A” on a test to defending my doctoral dissertation. While God made it work, it was not always easy. Despite moments when we were counting spare change for groceries, skipping celebrations to make ends meet, or managing the stress of life, she was my rock to get through it all. We were not always happy. We were not always in sync. We were, however, always together. She, truly, is my better half. It would be impossible, ungrateful, and unwise to pass this page up without acknowledging her. Thank you, my Love! (...--)

The lack of synchronicity tainted much of my work-life balance. Honestly, I did not know what I wanted to do when I got to graduate school. Additionally, I did not fit the mold either because I lacked traditional scientific curiosity. I would work on research because I saw it as a job, not as a passion. God's providence, however, gave me the mentors I needed: Ryan, Jeff, and Conner.

Ryan may not have always understood me, but he always supported me despite my non-traditional approaches to just about everything. It was through him that I was able to find my passion in graduate school: teaching. Because of my background, I had a perspective on how graduate education could be improved. Instead of telling me to toe the party line, Ryan told me it would be difficult. He told me it would take time. He told me to do it—though, probably through some obscure early 90's movie quote. The most important part of the mentorship is that while he was always my superior, I always felt like a peer. In some of my most heartbreaking moments in graduate school, Ryan did not judge. He was patient. He listened. He helped. I do not know if my actions will ever have a lasting impact on this department, but Ryan's actions have had a lasting impact on me.

I know I do not have the same patience as Jeff. I will never be able to bend the chopstick. However, his willingness to walk beside students in earnest support of their curiosity and understanding proved to me that graduate education was possible when done correctly. From the moment I asked him if I could do a mini lecture *in the class I was currently taking*, to backing my plea to be a GSI “one more time,” made me want to be an educator in this field. It may not have been a *Mr. Holland's Opus* moment when he finished his last lecture, but he truly did teach a master class on student-centered teaching. “Class dismissed.”

“It’s dangerous to go alone. Take this.” *The Legend of Zelda* starts with one door and then the rest of the world. Upon entering the door, a man gives Link a sword. It is that sword that starts the journey and the above quotation that I have associated with Conner. I did not have much time to work with Conner while in my undergraduate and I have no idea what he was thinking when he agreed to work with me. I did not even know what “bioinformatics” was then, but I knew that Conner did. I do not know if I could say that I actually did “bioinformatics” in my undergrad, but I did do my first “hello, world” and that legitimately changed my world. Conner has an honest fervor to see his students succeed; whether a publication, an alumni talk, or my defense, Conner celebrates those moments in such a genuine way that I continue to call him a mentor. Thank you, Conner.

I would also like to acknowledge my dissertation committee and the POISE Advisory Council (PAC). Unlike most students, I firmly fall within “alternative academic output.” My late game pivot away from traditional bioinformatics research toward educational research required my dissertation committee to challenge their own preconceptions of graduate student work. By no means was this an “easy ask.” I do not know whether Ryan was running some strong defense, but I never felt pressure or friction from my committee because of my decision. My presentations felt more like conversations. It seemed to me that the wise council that surrounded me was authentically interested in my work. This same feeling extended to the PAC. The PAC had no skin in the game. There was nothing in it for them. My success or failure would not impact their world in any real way. Nevertheless, they enthusiastically helped me shape POISE into what it became.

Furthermore, I want to acknowledge my editor, Katie Love. Through the support of the Pandemic Research Recovery program, I was provided the funding to assist in accelerating the

completion of my manuscript on threshold concepts and my dissertation. It was decided early on to earmark some of the funds to hire an external editor to review and revise this work from a purely grammatical and rhetoric standpoint. Her assistance made this process run so much more smoothly, and I couldn't be more appreciative. Thank you, Katie!

It would be prudent for me to acknowledge Dr. Vivian Cheung for showing me the characteristics of a mentor, scientist, and instructor that I would later use to shape what a proper mentor, scientist, and instructor should be. She will never know how her example shaped how I would approach all things academic from there on forward.

This work was funded in part by the Michigan Medicine Research. Innovation. Scholarship. Education (RISE) initiative at the University of Michigan and The University of Michigan Medical School Office of Research's Pandemic Research Recovery award. Raj, Paula, Nikki, and Helen...you took a gamble when you accepted a learner into the pilot cohort of RISE. I probably will never know exactly why you did that, and you will never know how much that enabled me become myself in graduate school. RISE afforded me the opportunity to earnestly pursue educational research within the medical school and catalyzed a profound change in both my dissertation as well as who I became as an educator. RISE tried to help us identify our champions, but you turned out to be some of my biggest. "Thank you" is not enough, but I don't know how else to say it. Thank you.

Lastly, to honor a wager, I would like to acknowledge DeLong & Co. Real Estate of Lansing, Michigan. At the same time as I was drafting this dissertation, my wife and I were getting our house ready to put on the market. We reached out to some local Dave Ramsey Trusted Agents and found Lara, Mike, and Arden. I told Laura that if she could get an offer on the plus side of our asking price before I defended that I would acknowledge them in my

dissertation. With absolute expediency and relative lack of stress, DeLong & Co. made good on their side of the agreement. As I do not bet often, this is one wager I was happy to “lose.” Thank you.

The actor Roberto Benigni (*Life is Beautiful*) once said, “It’s a sign of mediocrity when you demonstrate gratitude with moderation.” Therefore, without reservation, I express my genuine and enthusiastic gratitude: thank you for the journey.

Table of Contents

Dedication.....	ii
Acknowledgements.....	iii
List of Tables	xi
List of Figures.....	xii
List of Appendices	xiv
Abstract.....	xv
Chapter 1 Introduction	1
1.1 The current system	1
1.1.1 Teacher-centered education.....	1
1.1.2 “Teach the test”	5
1.1.3 Limited assessments.....	6
1.2 Three issues affecting bioinformatics education.....	9
1.2.1 The first issue: accessibility of content	9
1.2.2 The second issue: Effective curriculum design.....	13
1.2.3 The third issue: pedagogical content knowledge.....	16
1.2.4 Summary.....	19
Chapter 2 Bioinformatics Manipulatives.....	20
2.1 BAMnostic	22
2.1.1 The standards.....	22
2.1.2 The end-user experience.....	30
2.1.3 An OS-agnostic manipulative.....	31

2.1.4 The effects of BAMnostic	34
2.2 seqlogo	36
2.2.1 Sequence logo background.....	36
2.2.2 Dependency hell and programming cross compatibility	38
2.2.3 Sequence motif manipulative	40
2.2.4 The effects of seqlogo	41
2.3 Discussion	42
Chapter 3 Bioinformatics Interventions.....	43
3.1 Threshold concepts background.....	43
3.1.1 Features of threshold concepts	45
3.1.2 Threshold concepts in graduate education	47
3.1.3 Interdisciplinary science threshold concepts	49
3.2 Identification of threshold concepts within bioinformatics.....	49
3.2.1 Bioinformatics threshold concept study design.....	50
3.2.2 Bioinformatics student focus group.....	53
3.2.3 Bioinformatics student-centered survey	54
3.2.4 Bioinformatics faculty focus groups and survey	56
3.2.5 Limitations of data collection and analysis	57
3.3 Results	58
3.3.1 Student focus group.....	59
3.3.2 Student survey	62
3.3.3 Faculty feedback.....	69
3.4 Discussion	71
Chapter 4 Pedagogical Professional Development.....	74
4.1 Introduction	74

4.1.1 Faculty Survey.....	74
4.2 Pedagogy of Interdisciplinary Science Education.....	77
4.2.1 Theory of Change.....	78
4.2.2 Development and design	82
4.2.3 POISE Advisory council	85
4.2.4 Curriculum.....	85
4.2.5 Pilot cohort	88
4.2.6 Program evaluation.....	88
4.3 Results	89
4.4 Discussion	92
Chapter 5 Conclusion.....	96
5.1 Bioinformatics manipulative development	96
5.2 Threshold Concepts.....	98
5.3 POISE.....	102
5.4 Closing	102
Appendices.....	104
Bibliography	138

List of Tables

Table 1 Sequence Alignment Map (SAM) format.....	26
Table 2 BAM format.....	28
Table 3 Self-reported time in program from DCMB student survey.....	55
Table 4 Troublesome bioinformatics topics identified by student focus groups.....	60
Table 5 Faculty-refined list of troublesome topics provided by DCMB students.....	61
Table 6 Welch's 2-tailed t-test between junior and senior student responses.....	64

List of Figures

Figure 1 Action Priority Matrix:	4
Figure 2 Student operating system demographics.	12
Figure 3 Self-reported Student programming language usage (n=70).....	15
Figure 4 Cultural evolutionary model for pedagogical selection and transmission.....	18
Figure 5 Self-reported primary field of research of bioinformatics students (n=70).....	23
Figure 6 BAM indexing with Linear Index.	29
Figure 7 Basic BAMnostic output that demonstrates how the interface handles BAM files.	33
Figure 8 BAMnostic benchmarks	35
Figure 9 Example Position Probability Matrix (PPM).	37
Figure 10 Example sequence logo.	39
Figure 11 Bioinformatics threshold concepts study design	51
Figure 12 DCMB student feedback on statistical properties of bioinformatic data (n=56).....	65
Figure 13 DCMB student feedback on references to extant knowledge (n=63).....	66
Figure 14 DCMB student feedback on debugging strategies (n=54)	67
Figure 15 DCMB student feedback on sequence similarities of bioinformatic data (n=52)	68
Figure 16 DCMB-affiliated faculty survey of potential threshold concepts (n=13).....	70
Figure 17 DCMB-affiliated faculty survey demographic data (n=19)	76
Figure 18 POISE Theory of Change.....	81
Figure 19 POISE design components	84
Figure 20 POISE curriculum	86
Figure 21 Retrospective pre-post evaluation	90

Figure 22 POISE topic ranking.....	91
Figure 23 Threshold concept ontology	100
Figure 24 “Orphan” Concepts.....	101

List of Appendices

Appendix A Bioinformatics Curriculum Survey	104
Appendix B DCMB Threshold Concepts Faculty Survey	115
Appendix C POISE Pilot Cohort Feedback	133

Abstract

The education of bioinformatics, as an interdisciplinary science, can be negatively impacted by 1) incoming students attempting to negotiate complex concepts while lacking the educational scaffolding of the field's constituent disciplines while 2) taking breadth-first course requirements from 3) educators who may be lacking formal pedagogical training. Therefore, we attempted to ameliorate these issues through three different studies.

The first study was to identify common areas of code switching observed in required bioinformatics courses and develop software-based manipulatives to minimize the barriers to learning for students. The first manipulative we developed was BAMnostic, its purpose being to make genomic sequencing data used in bioinformatics courses and research more accessible by making the program easily installable across all major operating systems with no external dependencies. The second manipulative we designed was seqlogo. The purpose of seqlogo was to abstract complex software dependencies away from the students and allow them to focus on understanding and exploring sequence motif identification and analysis. As of this writing, both BAMnostic and seqlogo are required software for Department of Computational Medicine and Bioinformatics (DCMB) students at the University of Michigan and have been downloaded 210k and 20k times, respectively.

The purpose of the second study was to investigate bioinformatics curricula efficacy by first identifying potential bioinformatic-specific threshold concepts (if any), then suggesting curricular interventions and introducing pedagogical methodologies to address them. Threshold concepts (TC) are defined as troublesome knowledge that is transformative, irreversible, and domain-specific. Through a student-centered approach, the study began with student focus groups and surveys of students affiliated with DCMB to identify problematic concepts within the bioinformatic curriculum. These potential threshold concepts were then refined by direct collaboration with bioinformatics faculty. We received survey responses from 70 bioinformatics students (40% response rate). Students identified five conceptual bioinformatic obstacles: sequential data analysis; statistical distribution(s) identification and application; data ingest,

exploration, and management; data scaling; and references to extant knowledge. We collected 19 DCMB faculty survey responses (53% response rate) that suggested that while each identified concept was transformative, troublesome, and fundamental to understanding bioinformatics, none were bioinformatics specific. These findings corroborate other TC interdisciplinary science research suggesting that interdisciplinary fields may not have unique TCs.

The final portion of the research was focused solely on designing and developing the Pedagogy of Interdisciplinary Science Education (POISE) training program. The purpose of POISE was to address the gap in professional development specifically regarding the instruction of graduate-level students within the biomedical interdisciplinary sciences. This was a long-term approach to shift professional biases within the biomedical sciences towards a community of practice that supports and incentivizes pedagogical professional development in accordance with cultural evolution theory so that students who expect and respect well-trained educators at the graduate-level will themselves become well-trained educators to future students. This premise serves as a positive feedback loop that could potentially shift academic cultural norms and values. The POISE pilot cohort was comprised of 11 trainees and completed with a 100% completion rate. Trainees identified the most meaningful training as application of learning theories and authentic assessment and evaluation with virtual classroom management and technology in the classroom being the least. These findings and observations are currently planned to become embedded into the Medical Educators Novel Teaching On-Demand Resource (MENTOR) initiative in the University of Michigan Medical School.

Chapter 1 Introduction

Disclaimer: In no way is this dissertation meant to be a scathing indictment on the education that I have received at the University of Michigan or the Department of Computational Medicine and Bioinformatics. Quite the opposite. This dissertation was made possible only through the cooperation, support, and guidance of the faculty. Any critique or criticism is desired to be purely constructive and to ask, “how can we be better?”

For most graduate students, scientific higher education can be unequally divided into two components: research and coursework. Research output—in medical schools specifically—often dictates graduate student outcomes¹, thereby becoming the principle focus for graduate student education. Because each degree-granting program requires a specific number and type of courses to be completed towards the beginning of the graduate career, the unique unstructured learning environment and general lack of temporal guidance of the graduate program can become significant barriers to learning for students who come from more classical teaching methodologies². Therefore, it is important to explore the current system of graduate student education to diagnose potential points of failure within its pedagogical methodologies and identify educational strategies that may yield greater impact.

1.1 The current system

1.1.1 Teacher-centered education

Teacher-centered education is an approach that puts the teacher at the front of the classroom in an active role while the students observe and interact passively. This approach to

teaching has its merits, such as ease of classroom management, systematized sharing of knowledge, and general focus of content. Teacher-centered education is conventional Western education and has proliferated throughout the entire education hierarchy. This convention is also reflected in graduate education. Graduate courses are often taught by tenure-track faculty, who often have legally defined teaching loads³ and a conflicting emphasis on research⁴. As much of their effort is predisposed towards research, faculty tend to default to the known conventions of a teacher-centered style of teaching. The merits of this approach, however, no longer outweigh the drawbacks due the nature of graduate education.

Graduate education is predicated on students becoming “producers” of knowledge at some point. Therefore, many programs have core competencies⁵ with strong emphases on critical thinking in research and communication, with less emphasis on the knowledge of the field⁶. That is not to say that content knowledge is not important, but rather how one interacts with content knowledge is more impactful at this level of their education. This subtle shift in pedagogical focus and the unique makeup and number of the student body of graduate programs often negate the merits of teacher-centered education.

In education, “scaffolding” refers to a framework of pedagogical support logically constructed by an instructor that enables students to begin to process higher levels of learning within or around a subject. At the early stages of a graduate program, when students need scaffolding to adjust to the new expectations of graduate school, they generally have none. At the simplest level, graduate students are encouraged to communicate and collaborate effectively in their field, therefore their education should mirror these expectations and develop the proper scaffolding for the students to build upon. Therein lies an apparent dichotomy between the merits of teacher-centered education and the needs of graduate education: teacher-centered education

aids focus and classroom management by encouraging a passive student audience while graduate education requires the smaller, tightknit group of students to communicate, collaborate, and think critically at the cost of classroom management and potential loss of content coverage (e.g., breadth of curriculum). This single inconsistency in educational alignment heralds a greater concern not readily observed.

Possibly the most pernicious issue in graduate education is lack of instructor buy-in. As mentioned above, the faculty who often instruct graduate students tend to default to teacher-centered education due to a professional emphasis on academic output (i.e., “publish or perish”). While teacher-centered education necessitates a significant up-front investment of time and effort to develop the scope and sequence appropriate to the curriculum, it also enables educators to evaluate their action priority matrix^{7,8} (**Figure 1**) and determine their own level of involvement with student interaction. Additionally, this up-front cost is often a one-time expenditure with little revision between courses. Traditionally, this approach has empowered faculty to confidently determine whether the expected content has been covered and certain benchmarks have been met. Unlike all preceding levels of education with state/federal mandates, all graduate level benchmarks are self-imposed. Therefore, improper assessment and evaluation without oversight is another vestigial characteristic of teacher-centered education that plagues the effectiveness of graduate education.

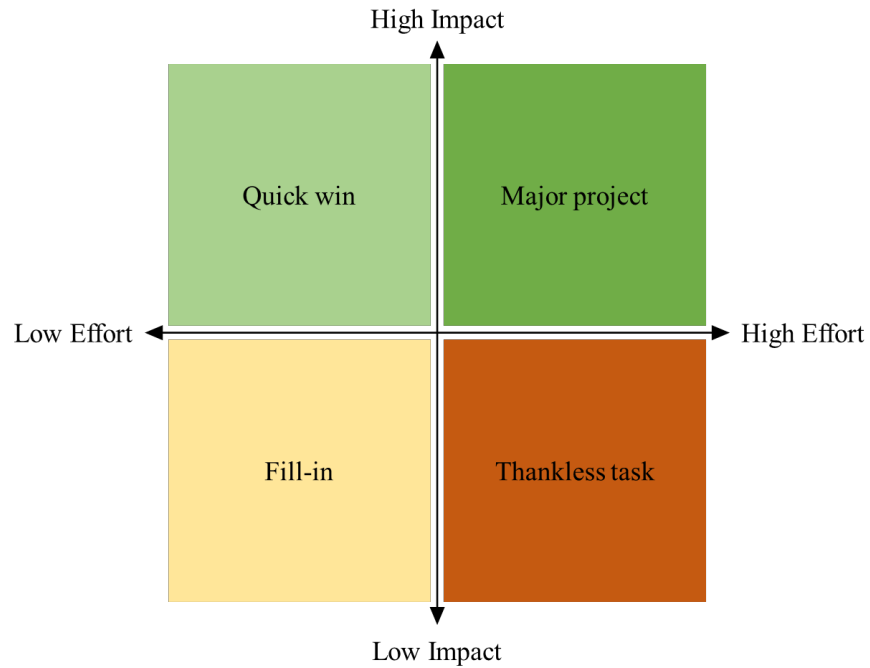


Figure 1 Action Priority Matrix:

A diagramming technique that allows a user to score a task by level of effort and degree of impact and choose tasks that are most efficient to their time. A low effort and low impact task are considered as “fill-in” or “busy work,” whereas a high effort and high impact task is considered as a “major project.”

1.1.2 “Teach to the test”

In 2002, the United States signed the No Child Left Behind (NCLB) Act into law. The purpose of NCLB was to attempt to address the education disparity between disadvantaged and high-performing students/schools by increasing accountability of schools via student outcomes by requiring any school receiving federal funding to administer standardized tests⁹. In 2015, the Every Student Succeeds Act (ESSA) was passed into law to shift the definition of student progress and the role of accountability from the federal government to the state. Additionally, ESSA requires schools to prepare students for college and careers by imposing a higher academic standard. To maintain accountability, states are still required to report the standardized testing results for both math and reading^{10,11}.

Through the influence of these laws, preprimary to postsecondary education is affected by standardized testing. Not only are student success and academic competitiveness measured by their scores but also the school and faculty’s successes are measured by these same student outcomes⁹⁻¹¹. Everything from a school’s funding to an instructor’s promotion, raise, or even position is weighted by the standardized test scores of its constituent students. Therefore, because of NCLB and ESSA, academia focused curricula to maximize student outcomes on standardized tests in a conceptual approach called “teach to the test¹².”

The premise of “teach the test” is that faculty are encouraged to concentrate on teaching students how to pass benchmarks instead of teaching the subject¹²⁻¹⁴. The reasoning of this concept is that a content standard (e.g., content covered in the tests) naturally tightens the curriculum—thereby reducing the faculty’s ability to create and explore subjects with their students¹². A side-effect of “teach the test” is that faculty merely teach students how to take tests.

A major hallmark of standardized tests is their modality: multiple-choice questions. To evaluate the validity of results, standardized tests tend to employ large banks of multiple-choice questions¹⁵, a form of question that can be designed to elicit higher order thinking or simply evoke rote memory¹⁶. These testing structures allow for rapid grading and objective scoring at the cost of student creativity and critical thinking. Moreover, if certain testing strategies are employed, it is possible to pass a multiple-choice test without knowing the answers¹⁷. Therefore, teaching how to test becomes a value proposition for both students and faculty alike: learn how to take the test and achieve better outcomes regardless of content mastery. Testing strategies, however, underpin a larger issue alluded to earlier: the tightening, or narrowing, of curriculum.

Anecdotally, while teaching at the graduate level the most exasperating question I regularly received from students is “Will this be on the test?”¹⁸ The reason this question frustrates instructors is that it highlights a fault in the system: emphasis on outcome in lieu of understanding the content¹⁹. The concept of grades in graduate school is different than that of any other educational level. Grades have less impact and less focus. Research output is weighed more heavily than grades²⁰ and the purpose of grades is more often used to check a proverbial box instead of assessing actual content proficiency. For example, a grade of a “B” is the lowest score a graduate student can receive to successfully fulfill course requirements²¹, yet grades in graduate school are curved such that approximately 90% of grades are at or above a “B.”²² Graduate school grade inflation has plagued higher education, reducing the actual meaning and importance of grades while also undermining the curriculum²³. Therein lies another apparent conflict in graduate education: students' hyper-focus on grades versus administrations' decreased focus on grades.

1.1.3 Limited assessments

The tension between disparate student and faculty views of grades in graduate programs is compounded by faculty continuing to assess and evaluate students using uninformative metrics—assessments requiring just basic knowledge or rote memory. Testing and assessment, when designed correctly, can be an invaluable resource in investigating student performance. At the graduate level, however, testing and assessment is rarely done correctly or appropriately^{24–26} starting with admission into a graduate program. As of this writing, one of the most hotly debated student outcomes is the Graduate Record Exam (GRE). The GRE has been a long held and heavily weighted metric in graduate student admissions. Recently, institutions have had to come to terms with the GRE and its effectiveness at determining potential student success regarding student admissions^{24,27,28}. This is an example of an uninformative student outcome: a good grade on the GRE does not necessarily mean the student would succeed in graduate school.

While the GRE is an uninformative gatekeeping metric for entrance to graduate education, within graduate programs, faculty, like their counterparts in lower levels of education, tend to default to teacher-centered education with institutions' encouragement of robust testing of its students for reporting purposes. As a corollary, faculty often employ testing strategies within their teacher-centered courses to gauge student understanding and progress in assigning student grades.

All assessment and evaluation should stem from teaching objectives of the course or lesson²⁹. Those teaching objectives leads to careful consideration of not only what is to be evaluated, but also whether the modality of evaluation appropriately measures the objective. Therefore, a course requiring higher order thinking should align its assessments and evaluations to ask higher order questions³⁰. Due to the demands of teacher-centered approaches and with the limited oversight on student testing within graduate courses²⁶, however, faculty can generate

assessment strategies not aligned with current teaching objectives. These strategies can result in handing a student a journal article and effectively asking them to do repeated document searches to answer a brief list of questions that do not evoke higher order critical thinking skills, making the assessment a moot exercise for both the student and faculty member alike. Additionally, these forms of assessment and evaluation do not capture the off-target effects of the unique nature of graduate education, like meaningful learning and tenacity^{25,26,31}.

As student-to-faculty ratios are often small in graduate school, the instructor-student interaction and dialogue have shown to have a significant impact on student outcomes and understanding^{32,33}. Furthermore, teacher-student interaction impacts student outcomes outside of test score-based metrics like attendance, tenacity, and grade progression³¹. Additionally, the pedagogical content knowledge (PCK) of the instructor has been investigated for its effectiveness at influencing student outcomes^{13,34-38}. These characteristics of impactful learning and education demonstrates another conflicting paradigm between teacher-centered education and graduate education: graduate student outcomes are influenced by how they interact with faculty as well as how proficient faculty is at instruction, whereas traditional teacher-centered education at the graduate level discourages teacher-student interaction and disincentivizes faculty pedagogical professional development.

Broadly speaking, the purpose of graduate education has been to transform consumers of knowledge into producers of knowledge. Therefore, assessments in graduate education ought to be fundamentally different than all previous stages of education as the expected result is different.

1.2 Three issues affecting bioinformatics education

In section 1.1, an apparent dissonance between the demands of graduate school educational output and the general application of classic pedagogical methodologies within a graduate educational setting was broadly detailed. In this section, it is posited that additional issues currently effect bioinformatics curriculum specifically at the graduate level. These are divided into three issues: 1) accessibility of content (“how students learn”), 2) effective curriculum design (“what students are learning”), and 3) pedagogical content knowledge (“who is teaching the students”).

1.2.1 The first issue: accessibility of content

The term “code switching” in education refers to the practice of alternating between two or more languages—or varieties of language—in a conversation³⁹. Code switching can be as simple as using a Western colloquialism to demonstrate a concept to a group of people who may or may not be native English speakers. Code switching can also extend to how one greets or gestures to people depending on their ethnicity. For example, a diplomat may bow to someone with an Asian background yet shake hands with an Anglo-Saxon. The effects of code switching in educational settings, however, have been found to have detrimental effects to student understanding⁴⁰⁻⁴², which relates to sociological concepts like stereotype threat and inclusivity³⁹.

Code switching is not inherently negative and can be used successfully to help develop language skills to learners in a controlled environment⁴⁰⁻⁴². Code switching, however, is both a conscious and subconscious act that can become a barrier to understanding for students when not addressed. For example, an instructor, to aid an explanation of a difficult concept like procedural writing, may ask students to write out the steps of making a peanut butter and jelly sandwich. For Western cultures, this is a straightforward exercise. However, for students from non-Western, it

is not as the learner is tasked with both comprehending what procedural writing is (like the other students) while also attempting to understand what a peanut butter and jelly sandwich is. This exercise now demonstrates an unintended and unequal learning burden based on culture.

Similarly, computational education and research are impacted by code switching because there exist multiple programming languages and diverse operating systems (OS). Therefore, learners can be classified by their level of proficiency with a given programming language and their familiarity with an OS. For example, teaching the programming language C++ to a student who uses a Windows-based laptop and has no previous experience in a programming language like C++ results in a higher educational burden than that of a student who uses a MacOS-based computer and has experience with another C-like declarative language. Both students are expected to fulfill the same requirements but with unequal effort.

The effect of code switching is magnified in graduate level interdisciplinary science programs. Traditionally, undergraduate education is generalized, and while there has been a rise in undergraduate interdisciplinary science programs like bioinformatics, there are far more students with generalized degrees⁴³. When these students are admitted into interdisciplinary science programs, they are often expected to overcome significant educational deficiencies. For example, a computer science student may understand the computational aspects of the research while having difficulty processing the biology and vice versa for a biology student. Bioinformatics, as a field, is characterized by the use of statistics, mathematics, information technology, and computer science to interrogate and answer biological questions. While the biological content and statistical/mathematical approaches vary from study to study, computational approaches are similar from student to student. Since computer science becomes

the academic burden for incoming students of bioinformatics⁴⁴⁻⁴⁸, it is one of the most prominent cases of code switching in early coursework.

In **Figure 2**, students across multiple bioinformatics courses self-reported (n=113) their primary OS. Of those respondents, 43% reported Windows as their primary OS. However, almost all bioinformatics toolkits, software, and computational resources are Linux-based^{49,50}. Coupled with generational differences in understanding computer architecture and file handling⁵¹, students in introductory or remedial bioinformatics computation courses are forced into a code-switching environment. From the moment they turn on their computers, they must relearn how a computer works just to attempt to understand the content of a course, posing a potentially significant barrier to learning for incoming bioinformatics students. Further exacerbating this inaccessibility of content is that there is little consistency of scaffolding: one course may be taught through the lens of one programming language while the next course is taught through another. The lack of intentional curriculum design and constant code switching ensures content coverage for the program at the cost of content internalization and integration for the student.

Which operating system does your computer use?
(n=113)

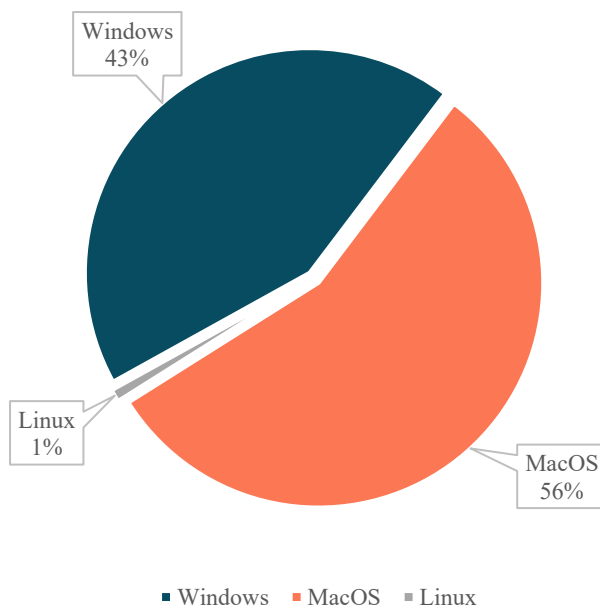


Figure 2 Student operating system demographics.

Pie chart represents data collected from voluntary entrance surveys of DCMB students across three different required bioinformatics courses (n=113) who self-reported the operating system they use on their computers. The distribution of MacOS (56%) and Windows (43%) demonstrate that while most high-performance and distributed computing solutions are Linux-like environments, a large portion of users personally work on a dissimilar operating system. The data of these surveys were collected for informational purposes only during the administration of the class and not as part of the larger work presented in this paper.

1.2.2 The second issue: Effective curriculum design

In section 1.2.1, the inaccessibility of content was explored through the lens of classroom code switching. To put it simply, section 1.2.1 explored “how students learn.” The important subsequent point is “what are students learning.” As an artifact of teacher-centered education and academic freedom⁵², little oversight of curriculum consistency and scaffolding is enforced at the graduate level^{53,54}. Students can often feel isolated due to the unique personalized—or ad hoc—training they require to align with their research interests⁵³. The resulting effects of this ad hoc education are programs with disparate courses with little overlap of scaffolding.

In any given domain, there are core concepts that are the foundation on which all subdomains are built⁵⁵⁻⁶³. These courses tend to chase local “aha moments” (moments of sudden insight or discovery) within their specific subdomain⁶⁴ instead of designing a holistic curriculum that addresses global—or fundamental—concepts. A curriculum that fosters more time and intentionality towards these concepts could potentially mitigate the cognitive load⁶⁵ on its constituent students as they begin to explore more specific hierarchical concepts later. Furthermore, student internalization and integration of these concepts is attenuated by the instructional consistency. In bioinformatics, for example, a student may be wrestling with the concepts surrounding sequential data generation (e.g., how a genome is sequenced) at the same time their instructor covers the use cases of the negative binomial distribution for RNA-seq analysis⁶⁶. Because of the split focus of both a fundamental concept of bioinformatics and the subdomain-specific concept, student understanding is likely attenuated^{57-59,62,64}.

Another aspect of the issue of effective curriculum design is the consistency of scaffolding, mentioned several times already. Since computer science is the backbone of

bioinformatics⁴⁴⁻⁴⁸, it bears examination that how bioinformatics is taught computationally matter just as much as what is taught. In **Figure 3**, data from a 2020 survey of the University of Michigan's Department of Computational Medicine and Bioinformatics students (n=70, **Appendix A**) shows that several programming languages are used by the student body. Of the identified programming languages, Shell/Bash, R, and Python represent the supermajority based on the responses. Therefore, it stands to reason that the courses taught within the department should support this supermajority by focusing its efforts on teaching the languages predominantly used. For the most part, this is true. As of this writing, however, many special topic courses (e.g., machine learning) within the same department use MATLAB or C++. This variety ensures that incoming students are already at a detriment despite taking the required introductory courses (see section 1.2.1).

Herein lies the issue of graduate education: effective curriculum design may unintentionally affect student outcomes and understanding. The onus of effort should be addressed. Is it the student's responsibility to overcome these educational gaps or should the educators make intentional and collaborative choices in how they design their curriculum with the students in mind? If it is the former, the status quo is acceptable. If it is the latter, then bioinformatics programs ought to design their curriculum with these concepts and drawbacks in mind.

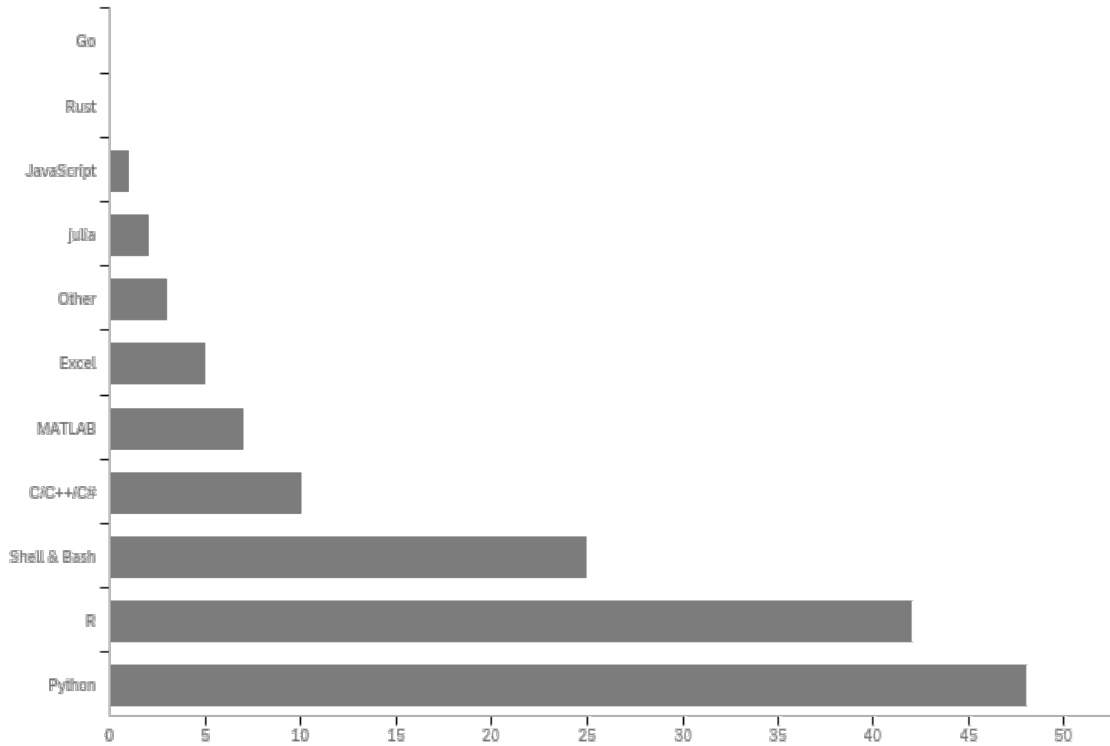


Figure 3 Self-reported Student programming language usage (n=70).

As part of a local student survey regarding bioinformatics educational effectiveness and curriculum design (**Appendix A**), students self-reported the programming languages they use in their research. The y-axis represents the language, and the x-axis represents the total count.

1.2.3 The third issue: pedagogical content knowledge

The final core issue of bioinformatics education is the pedagogical content knowledge of bioinformatics faculty. Pedagogical content knowledge (PCK) is the unique understanding and approaches a teacher uses to support a student's learning of a subject matter^{37,38,67}. In other words, PCK is not the subject matter to be taught, but rather how it is taught. It has long been understood that teaching practices, approaches, and student interactions positively impact student learning^{32,33,68}. Graduate-level bioinformatics faculty, however, are less likely to pursue professional development in pedagogical training due to cultural and institutional pressures despite mounting evidence that suggests its importance for student outcomes⁶⁹⁻⁷⁵.

This cavalier approach to graduate education results in teacher-centered faculty (see section 1.1.1) who produce courses that have little academic impact, which is an apparent contradiction to the faculty's own career progression, considering the presence of a teaching component for most tenure considerations. A faculty member who is evaluated in part via student interactions ought to concentrate on optimizing the efficacy of those interactions. This same approach, however, appears to be positively selected for in research-focused institutions.

When a cultural evolutionary model is applied to graduate education and career progression, poor or inappropriate pedagogical methodologies are likely to be perpetuated when pedagogical training is absent⁷⁶. The work of Grunspan, Kline, and Brownell (2018) describe a conceptual model (**Figure 4**) to discuss such a system, suggesting that graduate degree granting institutions tend to also be research-focused institutions. Upon graduation, potential faculty choose between teaching-focused and research-focused institutions. This first career transition introduces the premise of cultural evolution theory. At each major career transition (e.g., promotion), cultural and institutional pressures further select for specific traits in individuals;

that is, a research-focused institution will likely select faculty who are also research-focused.

This model also supports the concept that research-focused institutions incentivize research such that it discourages pedagogical professional development of its constituent faculty. Moreover, as graduates move to teaching-focused institutions, their influence on the graduate education system is no longer nearly as impactful since there are few bi-directional channels of communication between the two types of institutions⁷⁶. Using their model, Grunspan, Kline, and Brownell (2018) identified the most impactful places to introduce and incentivize pedagogical training: 1) PhD and postdoctoral students who are likely to become research-focused faculty, 2) faculty who train graduate students in research-focused institutions, and 3) the channels between teaching- and research-focused institutions⁷⁶.

Finally, recent political and cultural discourse has demonstrated that more nuanced and informed interpersonal approaches are a requirement to develop a more welcoming, diverse, equitable, and inclusive educational environment. Faculty can no longer just be subject matter experts but must also be trained how to manage their classrooms properly with these principles in mind. Trauma-informed teaching, classroom management, and inclusive teaching are necessary actions to manage tenuous situations and diverse demographics constructively⁷⁷.

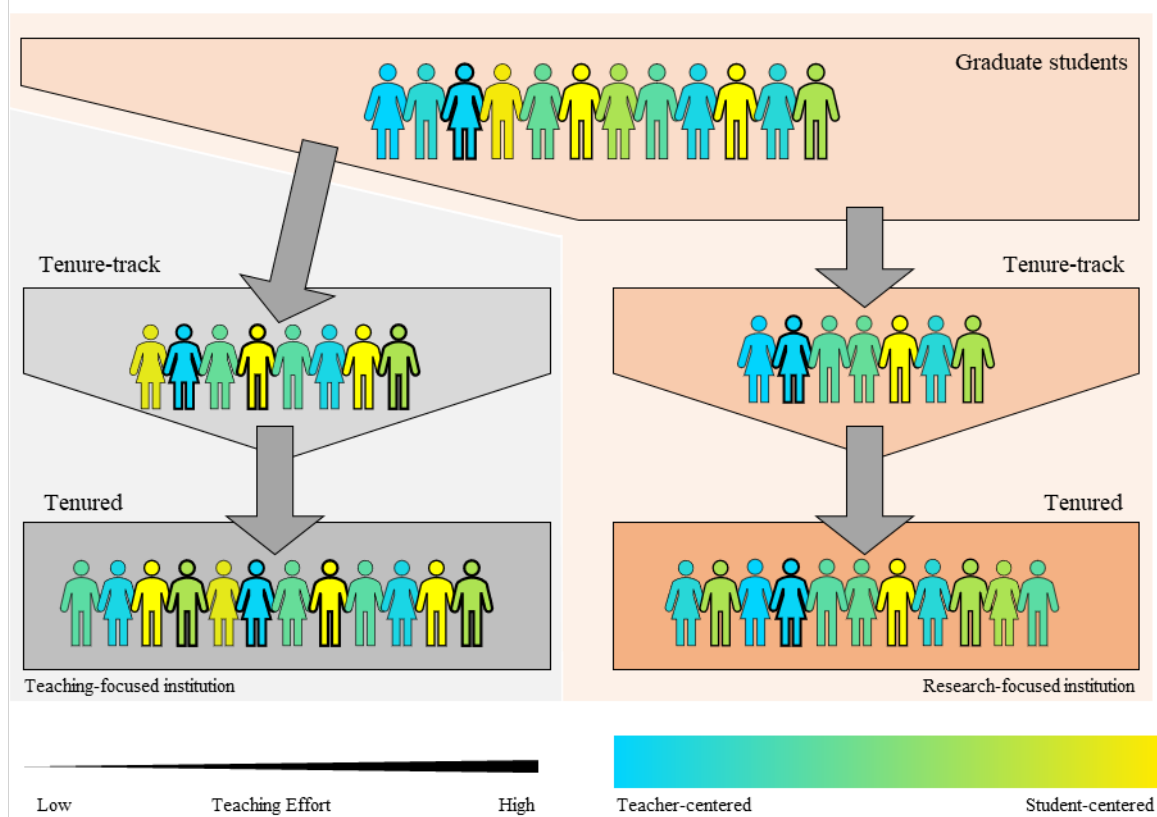


Figure 4 Cultural evolutionary model for pedagogical selection and transmission.

Modeled after the work of Grunspan, Kline, and Brownell (2018)⁷⁶ which defined a visual representation of a cultural evolutionary model of the spread pedagogical philosophies and practice within academia. The color of the individual represents their approach to teacher vs student centeredness and the width of the individual's border represents amount of teaching effort.

1.2.4 Summary

We attempt to address the three issues affecting bioinformatics education (section 1.2) in this dissertation through the research and development of three different interventions. To explore the issue of accessibility of content (section 1.2.1), we created multiple software-based manipulatives to ease technological barriers to entry and to mitigate pedagogical code switching. To investigate bioinformatics curricula efficacy (section 1.2.2), we conducted numerous student and faculty surveys to identify fundamental—but intractable —concepts that bioinformatics students are expected to internalize. Last, we designed a pedagogical training program for interdisciplinary scientists to address the lack of pedagogical content knowledge (section 1.2.3). We assert that intentional and explicit examination, evaluation, and integration of these interventions (or others like them) could result in transformative and systemic changes to the current bioinformatics graduate student education (section 1.1).

Chapter 2 Bioinformatics Manipulatives

This work, in part, was published in the Journal of Open Source Software 3(28) as I was the developer and first author of BAMnostic.

In education, a “manipulative” is a hands-on tool employed by instructors to supplement a lesson via constructivist-based active learning⁷⁸. Generically, a manipulative is something concrete used to teach something abstract. Mathematics and early childhood education regularly use manipulatives⁷⁹, using anything from clock dials and abacuses to develop understanding of abstract concepts like time and counting, respectively. Technology, however, has altered the form factor of manipulatives by enabling *in silico* applications. These “virtual manipulatives” shift the requirement of a manipulative’s being concrete. The premise of a manipulative is that, when utilized, it accelerates understanding of an abstract concept. As a corollary, manipulatives do not impact all students equally because students often conceptualize differently and internalize at different rates.

In addressing those student differences, manipulatives can also aid differentiation. “Differentiation” is the pedagogical technique of modifying instruction on an individual basis to meet multiple students’ needs⁸⁰⁻⁸². For example, if a student better understands a concept, one form of differentiation is to allow that student a chance to teach it to a fellow student or explore it on their own. If a different student, however, is having difficulty with the same concept, they may do additional exercises, participate in one-on-one instruction, or use a manipulative. Generally, differentiation in the classroom is achieved by tailoring either the content, process of learning, academic product, or learning environment⁸⁰⁻⁸². Differentiation can also be employed

to enhance inclusivity of students and accessibility of content in the classroom by enabling instructors to address complex concepts and diverse learners on a case-by-case basis. These approaches can make learning more equitable⁸³.

Graduate-level interdisciplinary science education is inherently different from previous educational modalities². Unlike more basic fields of scientific research (e.g., biology or physics), in which graduate learning becomes a more organic extension of previous work, interdisciplinary science often requires students to overcome potentially significant educational shortcoming due to requirements in less familiar disciplines⁴³. Of the many requirements in interdisciplinary science, computer science is the keystone of the computational fields, scilicet bioinformatics⁴⁶. This is because although both the biological content and statistical/mathematical approaches may differ between studies, the computational methodologies are largely conserved⁴⁴⁻⁴⁶. Therefore, the most impactful concepts for incoming bioinformatics students (e.g., high-throughput data access and analysis) are likely to be centered on computer science creating the need for differentiated instruction for students with diverse backgrounds.

This chapter focuses on our approaches to mitigate the cognitive load⁶⁵ of these students by developing pedagogical tools that aid in understanding abstract computational concepts^{79,84,85}. Because of the nature of bioinformatics graduate education, these tools operate not only as contemporary analysis programs for students' given disciplines, but also—and more importantly—as pedagogical manipulatives designed specifically to bridge complex conceptual gaps. Since conceptual code switching and technological diversity between students may lead to unequal educational onboarding⁴⁰⁻⁴² (section 1.2.1), we created two manipulatives with the express purpose of facilitating learning by abstracting away confounding technological barriers to entry. These two tools are named “BAMnostic” and “seqlogo” and were designed specifically

for bioinformatics graduate students but have also seen additional application in the general bioinformatics research community.

2.1 BAMnostic

The field of bioinformatics is divided into many sub-fields like genomics (the study of structure, function, and mapping of genomes)⁸⁶, epigenetics (the study of changes in organisms caused by modification of gene expression)⁸⁷, transcriptomics (the study of the entire set of transcripts expressed by a cell, tissue, or organism)⁸⁸, etc. According to a 2021 student-wide survey of students within the Department of Computational Medicine and Bioinformatics at the University of Michigan (n=70; **Appendix A**), many of the students conduct research largely within Genomics (26%), Computational Biology (19%), and other fields closely associated with genetics (**Figure 5**), suggesting that most students will likely encounter sequencing data within both their course work as well as research.

2.1.1 The standards

Most sequencing data are biological data relating to the genes, transcripts, and proteins that have been coded to allow for computational analysis. Unlike a dictionary or a large-scale survey in which the individual entries or questions are logically unrelated to those flanking them, the order of the sequencing data matters. In other words, the data relating to the left most position of a gene should be near—but before—the data from the right most position of the same gene. These data are produced by several technologies like short-read (e.g., Illumina, 150-300 base pairs) and long read (e.g., Oxford Nanopore, up to 100k base pairs)⁸⁹ next generation sequencing platforms. The output of these platforms come in the form of “reads.” A “read” is an alphabetic sequence of biological data representing DNA, RNA, or amino acids. Downstream

analyses use a reference for a given sample to align the reads, thereby providing positional data for the reads⁸⁹⁻⁹³. As an analogy, alignment is like repeated—but complex—word searches in which the reads are the "words" and the reference being used is the "letter bank." Read alignments provide positional data and quality scores for all reads produced by the sequencing platforms—that is, we know what the read is (alphabetic sequence), where it belongs (numeric position), and how well it fits (quality score) relative to the reference.

Self-reported primary field of research of bioinformatics students

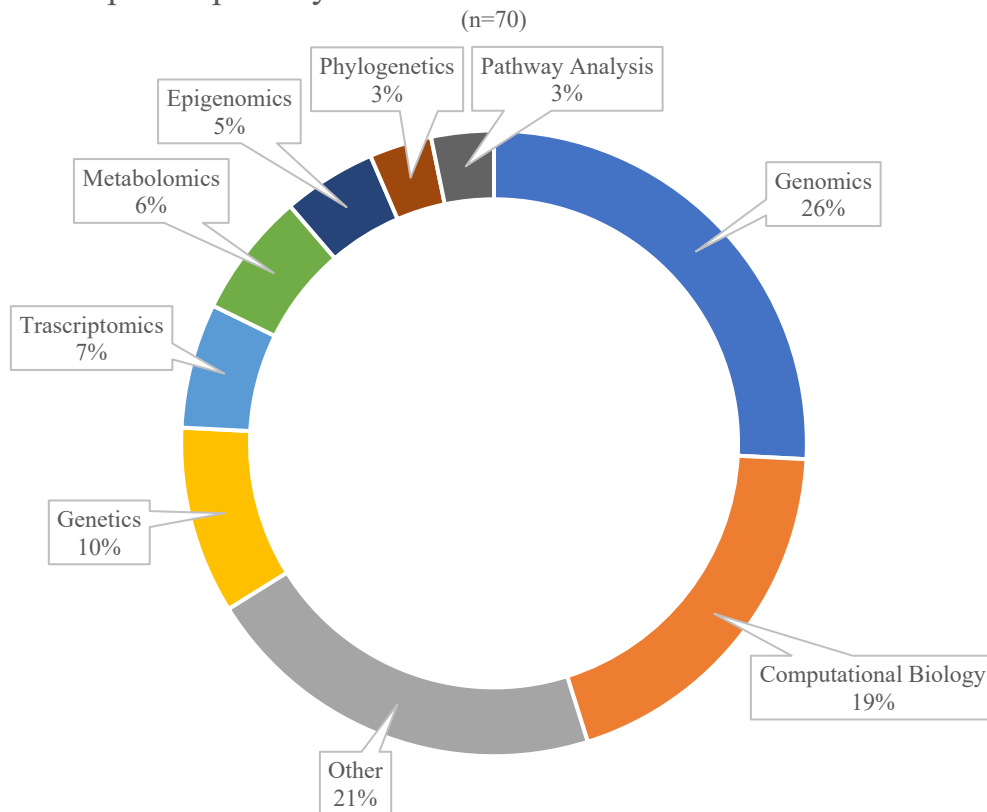


Figure 5 Self-reported primary field of research of bioinformatics students (n=70).

In a student-wide survey of all students associated with the Department of Computational Medicine and Bioinformatics at the University of Michigan in late 2021 (**Appendix A**), students were asked to self-report to the following prompt: “What primary sub-field of bioinformatics is either your current or prospective research in?” Genomics (26%) and Computational Biology (19%) were the highest reported of all the categories, while “Other” (21%) aggregates multiple singletons.

While the data can represent different types of analyses (e.g., DNA, RNA, methylation, etc.) the underlying data are similar enough that standard formats were developed. These standards are called FASTA⁹⁰ and Sequence Alignment Map (SAM) format⁹³. FASTA format is simply the name, metadata, and alphabetic sequence of a read as produced by sequencing platforms. An extension of this format is called FASTQ, where “Q” means “quality,” signifying that the quality score of the sequence is also captured⁹⁴. Whereas FASTA/FASTQ formats are the output of sequencing platforms, SAM format is used by sequence aligners and captures additional details about the read like its position relative to the reference and its potential relation to other reads (**Table 1**).

Because early sequencing studies produced much smaller datasets than contemporary research, both the FASTA/FASTQ and SAM formats were plain text, meaning the files created could be opened on computers with no special software and the data were human-readable. However, when considering that the SAM format further extends the information of a single read and that more efficient technologies led to higher throughput data were created, two major drawbacks of the standards were discovered: speed and size. Speed here represents the ability to quickly and accurately access a desired section of reads and was a drawback because many studies were only concerned with a relatively small portion of a given dataset at any given time. The second drawback of size is readily apparent since the number and size of studies that require sequencing implies an increase in space to store these new data.

#	Name	Description	Example
1	QNAME	Query template name	B7_591:8:4:841:340
2	FLAG	bitwise flags	73
3	RNAME	Reference sequence name	chr2
4	POS	1-based leftmost mapped position	1
5	MAPQ	Mapping quality	99
6	CIGAR	Compact idiosyncratic gapped alignment report string	36M
7	RNEXT	Reference name of mate/next read	*
8	PNEXT	Position of the mate/next read	0
9	TLEN	Observed template length	0
10	SEQ	Sequence of aligned segment	TTCAAATGAACTTCTGTAATTGAAAAATTCATTAA
11	QUAL	ASCII string of offset Phred-scaled base quality scores	<<<<<<<<;<<<<<<<<;<<<<<;<;:<<<<<<<<;;
12	INFO	Tag:type:value information	MF:C:18 Aq:C:77 NM:C:0 UQ:C:0 H0:C:1

Table 1 Sequence Alignment Map (SAM) format.

The SAM format details 11 mandatory entries for each linear alignment of a segment⁹³. There is an optional column of data that can take any other additional data that is not captured by the previous 11 columns.

To address these two major drawbacks, the Binary Alignment Map (BAM) format was developed⁹³. BAM format compresses the SAM data in two ways: binary compression of data and BGZF compression. The major tradeoff of this format is that the data are no longer readily human-readable. Binary compression converts plain text data into binary representations of the data (**Table 2**). BGZF compression is a special use case of the gzip file format⁹⁵, that is, instead of compressing the entire file using gzip, the file is broken up into equal (or smaller) blocks and then each block is gzip compressed individually. These blocks are then concatenated, thus forming the BGZF file format. Not only is the BGZF format a data compression technique, but it also, since the blocks are of a known maximum size and are concatenated onto each other, makes random access achievable. Random access is the ability to quickly seek a specific section of a file instead of starting at the top of a file and serially processing until the section of interest is found. Through the BAM format, both the drawbacks of size of sequencing data ($\approx 50\text{-}80\%$ compression over original⁹⁶) and the speed of accessing data (through random access; **Figure 6**) are overcome.

<i>Field</i>	<i>Description</i>	<i>Type</i>
<i>magic</i>	BAM magic string	char[4]
<i>l_text</i>	Length of header text	uint32_t
<i>text</i>	Plain header text in SAM	char[l_text]
<i>n_ref</i>	# of reference sequences	uint32_t
<i>List of reference information</i>		
<i>l_name</i>	Length of reference name + 1	uint32_t
<i>name</i>	Reference sequence name	char[l_name]
<i>l_ref</i>	Length of the reference sequence	uint32_t
<i>List of alignments (until end of file)</i>		
<i>block_size</i>	Total length of the alignment record	uint32_t
<i>refID</i>	Reference sequence ID	int32_t
<i>pos</i>	0-based leftmost coordinate	int32_t
<i>l_read_name</i>	Length of read_name below	uint8_t
<i>mapq</i>	Mapping quality	uint8_t
<i>bin</i>	BAI index bin	uint16_t
<i>n_cigar_op</i>	Number of operations in CIGAR	uint16_t
<i>flag</i>	Bitwise flags	uint16_t
<i>l_seq</i>	Length of seq	uint32_t
<i>next_refID</i>	Ref-ID of the next segment	int32_t
<i>next_pos</i>	0-based leftmost position of the next segment	int32_t
<i>tlen</i>	Template length	int32_t
<i>read_name</i>	Read name	char[l_read_name]
<i>cigar</i>	Compact Idiosyncratic Gapped Alignment Report	uint32_t[n_cigar_op]
<i>seq</i>	4-bit encoded read	uint8_t[(l_seq+1)/2]
<i>qual</i>	Phred-scaled base qualities	char[l_seq]
<i>List of auxiliary data (until end of alignment block)</i>		
<i>tag</i>	Two-character tag	char[2]
<i>val_type</i>	Value type	char
<i>value</i>	Tag value	(by val_type)

Table 2 BAM format.

This table represents the BAM format as described by Li et al⁹³. A BAM file is divided into three parts: the metadata, list of references, and the alignments (or reads). Any read can contain additional information within the auxiliary section so long as it follows the tag:value format.

BAM indexing with Linear Index

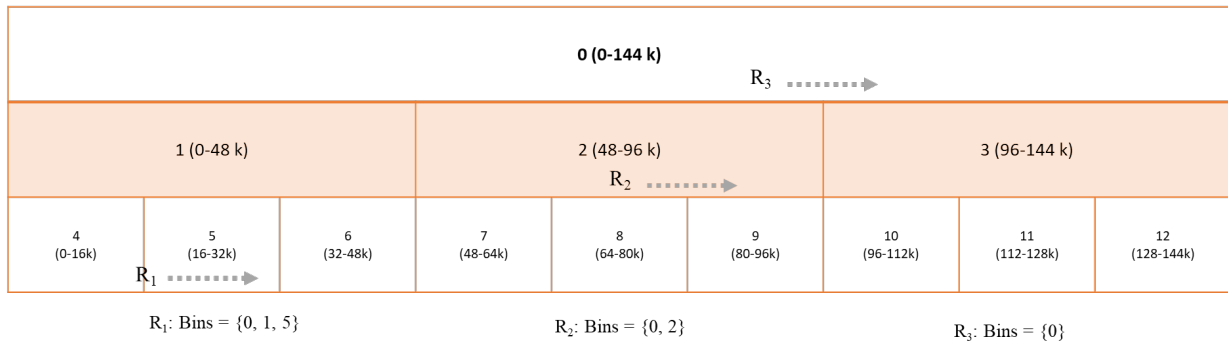


Figure 6 BAM indexing with Linear Index.

BAM indexing creates an accompanying file called the BAM Index file (BAI). The purpose of the BAI file is to capture the representation of the R-tree binning scheme and record the linear index. The binning scheme assigns bins to regions spanning 229, 226, 223, 220, 217, 214 base pairs. Smaller bins are subsets to larger bins. Each bin is also recorded with the smallest file offset of the alignments that overlap with the spanning region. Since the bins span base pair regions and contain the smallest file offset, a single seek call needs to be made to randomly access the file exactly at the region of interest. For example, since region of interest R_1 is contained within Bin 5 (region 16kbp- 32kbp), a single seek call to the file offset associated with Bin 5 is needed.

2.1.2 *The end-user experience*

As the SAM/BAM file formats became the *de facto* standards of sequencing data, tools were developed to make accessing and analyzing the data more straightforward. At the time of this writing, the top three packages downloaded on Bioconda (a channel for the conda package manager that specializes in bioinformatics software) are samtools, htlib, and pysam^{49,93,97}.

The samtools software was created to perform several low-level operations on SAM/BAM files like converting between formats, sorting and merging files, querying sections of the file or for specific reads, and showing alignments in text-based viewer⁹³. This software has become a linchpin in bioinformatics and genetics research for nearly as long as the SAM format has been around.

Htlib was born out of the samtools project as it was developed to create a library of tools that enable developers and scientists to leverage robust, performant, and standardized approaches to interrogating sequencing data⁹⁷ in much the same way that samtools does. These tools include decompression of data, random access, etc. Due to its utility, htlib is a major dependency for many other analysis toolkits, among them pysam.

Pysam was written to create an interface for Python programming language users by converting the low-level language used by htlib⁹⁸ into a more accessible language. Since Python is one of the primary programming languages used in bioinformatics research, pysam enables researchers to dynamically interrogate sequencing data with the ease of use that Python provides.

When considering the landscape of tools and software available for bioinformatics research, it is easy to overlook a crucial component of the user experience as it relates to bioinformatics pedagogy: infrastructure. These—and most other bioinformatics tools—are built to operate within HPC environments. These environments are mostly Linux-based. Therefore,

these tools inherently employ Linux defaults such as file structures, commands, and general lack of graphic user interfaces. However, as shown earlier (**Figure 2**), many students (43%) enter bioinformatics programs with only rudimentary understanding of non-Linux environments, like Windows. The simplest barrier between the two OSs is the file structure: Windows uses multi-forested drive volumes (e.g., C:\) whereas Linux uses single tree mounts (e.g., /mnt/c/). Additionally, as apparent by the previous examples, the two systems use different forms of “pathing.” Pathing is how the location of a file or folder is represented by the system. Windows uses backward slashes (“\”) to delimit paths while Linux uses forward slashes (“/”). At the most basic level, a Windows user experience with bioinformatics tools is already subject to the most straightforward version of code switching. Furthermore, many of the most common bioinformatics tools are incapable of being installed in native Windows environments, which means that while being the most downloaded packages from Bioconda, they cannot be installed on Windows systems. This pernicious cognitive dissonance and potential impediment to incoming students is the basis of why we developed our first bioinformatics manipulative.

2.1.3 An OS-agnostic manipulative

In developing lessons plans, one should always start from the end and work backwards, a system known as “backward design²⁹.” Backward design frameworks suggest first identifying the desired outcome of a lesson, then designing the means to observe/evaluate the presence or progress of that outcome, and finally developing the content to make such an outcome possible. Therefore, it is important to ascertain the final student outcome and develop lessons and manipulatives that support the achievement of that outcome. As detailed earlier, computer science is the common denominator for bioinformatics students, followed closely by use of programming language (**Figure 3**). Additionally, many students in bioinformatics programs are

likely to interact with sequencing data at some point. If a backward design of a bioinformatics program identifies these areas as part of the desired outcome, then course work must help students achieve that outcome

Therefore, to ensure that students have a better understanding of sequencing technology and its data, we developed a manipulative that would abstract away as many discrepancies in OS choice as possible from the students through a more accessible manipulative called BAMnostic. BAMnostic would be able to perform the most common operations that are employed by more mainstream tools like pysam. The key difference, however, is not that BAMnostic would be more performant than these other tools, but that it can be readily and simply installed across a wide range of programming environments. Last, BAMnostic would be written completely in Python for transparency and extensibility by students as most students currently wield Python as a programming language (**Figure 3**).

BAMnostic was written to be a fully featured, pure Python implementation of BAM file random access and parsing⁹⁹. Since pysam—the software BAMnostic was modeled after—could not be installed on Windows systems because it depends on htlib, BAMnostic was developed to have no dependencies: it can be installed as a standalone package with no other requirements. Consequently, BAMnostic can be installed on any system that can install Python, thus expanding the realm of sequencing data exploration from HPC environments all the way to smart phones and preventing students from experiencing dependency issues during initial classroom preparation. Last, to overcome the large amounts of legacy code in academic research, BAMnostic was written to support any Python version greater or equal to 2.7. As such, BAMnostic potentially makes genomic research and analytics available to a much greater software demographic⁹⁹. **Figure 7** demonstrates a simple example of BAMnostic usage.

```

1 >>> bam = bamnostic.AlignmentFile(bamnostic.example_path, 'rb')
2 >>> for i, read in enumerate(bam.fetch('chr2', 1, 100)):
3 ...     if i >= 3:
4 ...         break
5 ...     print(read)

6 B7_591:8:4:841:340 73 chr2 1 99 36M * 0 0
  TTCAAATGAACTTCTGTAATTGAAAAATTCATTTAA
  <<<<<<<<<;<<<<<<<<<;<<<<<<;<;<<<<<<<<;;
  MF:C:18 Aq:C:77 NM:C:0 UQ:C:0 H0:C:1 H1:C:0
7 EAS54_67:4:142:943:582 73 chr2 1 99 35M * 0 0
  TTCAAATGAACTTCTGTAATTGAAAAATTCATTTA
  <<<<<<<<;<<<<<<<<:<<<;<<<<<;<<<<;<<<<:<;<<<<5
  MF:C:18 Aq:C:41 NM:C:0 UQ:C:0 H0:C:1 H1:C:0
8 EAS54_67:6:43:859:229 153 chr2 1 66 35M * 0 0
  TTCAAATGAACTTCTGTAATTGAAAAATTCATTTA
  +37<=<.;<<7.;77<5<<0<<<<;<<<<27<<<<<<<<
  MF:C:32 Aq:C:0 NM:C:0 UQ:C:0 H0:C:1 H1:C:0

```

Figure 7 Basic BAMnostic output that demonstrates how the interface handles BAM files.

The grey numbers on the left represent line number for demonstration only and are not present in actual use. Line 1 connects the program to a BAM file. Lines 2-5 shows a simple random access approach to produce output of the first three coordinate-sorted reads present in the BAM file at a given location within a genome (e.g., chromosome 2 ('chr2') between base pair positions 1 and 100). Lines 6-8 are the SAM-formatted reads from the BAM file from the code above.

2.1.4 The effects of BAMnostic

As of this writing, BAMnostic has over 210,000 downloads and 75 stars on its GitHub repository (<https://github.com/betteridiot/bamnostic>)—the same amount as the original FASTA package⁹⁰ (<https://github.com/wrpearson/fasta36>). Additionally, BAMnostic has become a required install for students in the Department of Computational Medicine and Bioinformatics at the University of Michigan to specifically combat difficulties in technological on-boarding and has also become a dependency of numerous packages and repositories outside of its intended educational purpose.

Several off-target effects of the pure Python nature of BAMnostic are that 1) it allows BAMnostic to be readily adapted to be parallelized using standard Python libraries like multiprocessing and threading and 2) it allows BAMnostic to be implemented in PyPy—a Python implementation of Python. By leveraging the just-in-time compiler of PyPy, BAMnostic can observe noticeable speedups. These same benefits cannot be observed by pysam as pysam cannot be installed in a PyPy environment or make use of standard multiprocessing and threading Python libraries (**Figure 8**). This just shows, however, the extensibility and accessibility of BAMnostic but makes no claim at outperforming pysam as pysam makes use of C extensions and wrappers to achieve lower-level efficiency.

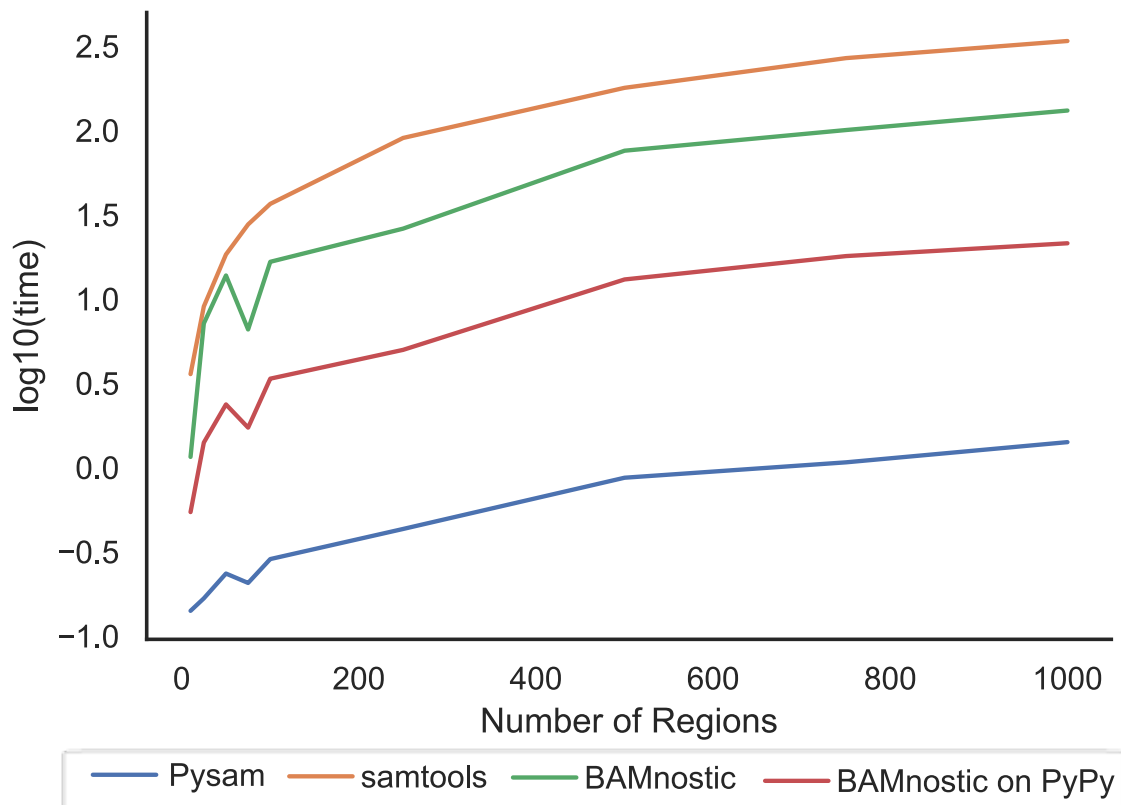


Figure 8 BAMnostic benchmarks

To benchmark BAMnostic, we used Ribosomal profiling data from GM19257 (SRR1585557)¹¹³ and performed multiple rounds random access sampling of n number of randomly selected gene expression BED regions. All times are scaled using \log_{10} . BAMnostic on Windows is not shown. These data demonstrate that while BAMnostic is not as performant as pysam, it is comparable and more versatile.

2.2 seqlogo

Along a similar vein as BAMnostic, the second bioinformatics manipulative that we developed is called “seqlogo.” This manipulative was designed specifically for the students of the BIOINF529: Bioinformatics Concepts and Algorithms course within the Department of Computational Medicine and Bioinformatics at the University of Michigan. BIOINF529 is a required course for all bioinformatics students and during the module of instruction on how to programmatically identify sequence motifs, we observed numerous students struggling with syntax, usage, and installation of existing software and packages—again stemming from differences in OS and Python versioning.

2.2.1 *Sequence logo background*

A sequence motif is a pattern present in nucleotide and amino acid sequencing data, and identifying motifs is a fundamental task in determining everything from protein structure to transcription factor binding sites^{100–104}. A sequence motif is identified using position-specific scoring matrices (PSSM). The most used PSSMs are 1) position frequency matrix (PFM; tallies each observation at each position), 2) position probability matrix (PPM; calculates the probability of each observation at each position given the PFM), and 3) position weight matrix (PWM; calculates the log likelihoods of each observation at each position). An example PPM of a 10 base pair length DNA sequence can be seen in **Figure 9**.

	A	C	G	T
0	0.16	0.10	0.06	0.68
1	0.02	0.15	0.81	0.02
2	0.21	0.06	0.08	0.66
3	0.15	0.22	0.07	0.57
4	0.82	0.05	0.08	0.05
5	0.44	0.02	0.05	0.49
6	0.17	0.05	0.49	0.29
7	0.20	0.07	0.19	0.55
8	0.18	0.07	0.42	0.33
9	0.28	0.49	0.20	0.03
10	0.14	0.06	0.02	0.78

Example of a PPM

Figure 9 Example Position Probability Matrix (PPM).

This example demonstrates what a possible PPM would look like for a 10 base pair DNA sequence.

Once a segment's PSSM has been computed, it is then graphically rendered into what is called a "sequence logo," which uses the counts/probabilities/weights for each observation (e.g., A, C, G, or T) at a given position and renders the observation's height relative to its score (**Figure 10**)¹⁰⁵. These sequence logos provide a fuller representation of motif than complex, text-based consensus sequences.

2.2.2 Dependency hell and programming cross compatibility

Until recently, the simplest method to generate visually compelling sequence logos was either through web-based interfaces (e.g., WebLogo¹⁰⁵) or through disparate third-party packages (e.g., seqlogo¹⁰⁶ in Bioconductor¹⁰⁷ of the R programming language). During curriculum development of BIOINF529, we recognized that the sequence motif module was going to suffer since the class was to be taught via the Python programming language. At the time, it required multiple third-party libraries (e.g., Numpy, pandas, biopython, MATPLOTLIB, etc.) to generate a single sequence logo. While many of these libraries are widely used, it necessitated in-depth knowledge of lesser used and/or maintained features of these libraries. In some cases, any attempt to generate a sequence logo would dissolve into a "dependency hell," the dilemma in which one package may interfere with the dependencies of a different one¹⁰⁸.

A cost/benefit analysis suggested that it was not worth the pedagogical investment to explain and troubleshoot individual computers for a one-off exercise. Furthermore, it would become a cognitive burden to the students to ask them to code switch between Python, R, and web-based platforms. Therefore, we decided to ameliorate the code switching and cognitive load imposed on our students by a creating Python sequence logo plotting program: seqlogo.

Plot the sequence logo with information content scaling

```
# Setting seed for demonstration purposes
>>> np.random.seed(42)

# Making a fake PWM
>>> random_ppm = np.random.dirichlet(np.ones(4), size=6)
>>> ppm = seqlogo.Ppm(random_ppm)
>>> seqlogo.seqlogo(ppm, ic_scale = False, format = 'svg', size = 'medium')
```

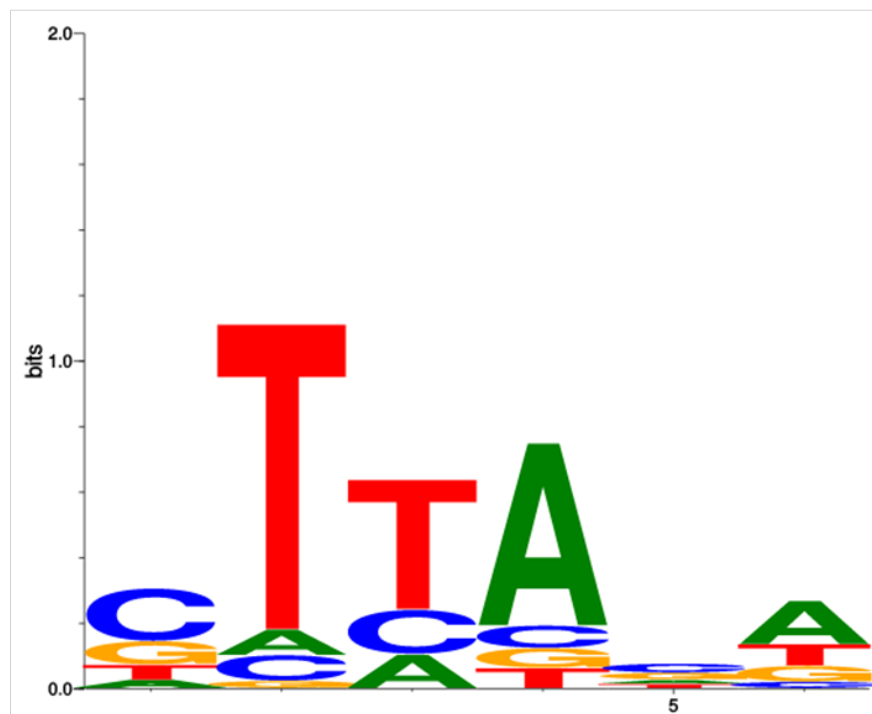


Figure 10 Example sequence logo.

The Python code demonstrates how to generate a random Position Probability Matrix and use seqlogo to format and render the data. A sequence logo will plot the observations with their height relative to their score in the matrix at a given position across all positions of the sequence. For example, at position 2, “T” is observed far more than any other observation.

2.2.3 *Sequence motif manipulative*

We intentionally developed our second manipulative, seqlogo, (<https://github.com/betteridiot/seqlogo>) to assist the students of BIOINF529: Bioinformatics Concepts and Algorithms at the University of Michigan in the Department of Computational Medicine and Bioinformatics. Seqlogo attempts to blend the user-friendly interface of the similarly named R package, seqLogo¹⁰⁶, with the powerful web-based API of WebLogo¹⁰⁵. We even branded seqlogo with the version number of 5.29.# to signify its original purpose of supporting BIOINF529. Seqlogo was listed on the Bioconda⁴⁹ channel for the conda package manager that many Python users are familiar with. This listing enabled our students to install seqlogo quickly with a single command while simultaneously allowing the conda package manager to manage package dependencies seamlessly.

Early iterations of seqlogo only rendered sequence logo plots (**Figure 10**) and nothing else, which consequently added to students' cognitive load as they had difficulty ascertaining which PSSM to provide as input and how it should be formatted. Ergo, we extended the design of seqlogo to include features that supported all types of PSSMs and the requisite PSSM conversion functions. This functionality gives seqlogo the flexibility to handle most standard sequence alphabets (e.g., DNA, RNA, and amino acid) as well as their reduced and ambiguous versions. Additionally, by piggybacking onto WebLogo's API, seqlogo can produce sequence logo plots as scalable vector graphics (SVG), portable document format (PDF), portable network graphic (PNG), and more. Seqlogo depends on a handful of third-party—albeit commonly used—packages: Numpy¹⁰⁹, pandas^{110,111}, and WebLogo¹⁰⁵. Two additional Linux-based

applications can be installed to further extend sequence logo rendering support on Linux systems: `ghostscript` and `pdf2svg`.

Figure 10 demonstrates a straightforward example of how to produce a sequence logo using `seqlogo`. The first two lines of code are for demonstration purposes to produce a PPM quickly such that the sum across all probabilities for a given position equals one. So long as the input array matches the dimensions of the expected sequence alphabet, `seqlogo` appropriately handles those inputs. Furthermore, `seqlogo` exposes a new entry point called a “Complete Position Matrix,” or “Cpm.” If the student submits any PSSM type to the Cpm, `seqlogo` will calculate all other forms of PSSMs and provide them within the Cpm object.

2.2.4 The effects of seqlogo

The importance of `seqlogo` is not in its ability to create sequence logos. Many other programs already do that. `Seqlogo`’s importance is also not that `seqlogo` makes sequence logos in Python. Recently, a library called `Logomaker`¹¹² was developed that is more robust and flexible in how it renders sequence logos in Python. The importance of `seqlogo` is that it was designed specifically as a manipulative to BIOINF529 students to offload the cognitive effort of getting a computer program to work and streamline the learning process of the concepts of sequence motifs and logos themselves. This marked a milestone in pedagogical innovation within graduate-level bioinformatics: the creation of a tool to aid in the instruction of bioinformatics instead of just the research of bioinformatics.

As of this writing, `seqlogo` has been downloaded over 20,000 times and has become a dependency within toolchains that extend outside of BIOINF529. For example, `seqlogo` is a dependency for both the Sequence Motif Enrichment and Genome Annotation Library (SMEAGOL)¹¹³ of Gruber Science Lab out of Germany as well as the Regulatory Sequence

Analysis Tools web resource¹¹⁴. This usage beyond BIOINF529 demonstrates that a tool like seqlogo, written for the express pedagogical purpose of ease of use and understanding, benefits those outside academia, the ultimate goal of research.

2.3 Discussion

This chapter discussed the concept, development, and usage of pedagogical manipulatives within the context of graduate-level bioinformatics. The purpose was to exemplify the use cases of how to make impactful tools that support graduate-level learning while reinforcing real-world bioinformatics research. BAMnostic was written to make genomic data stored in BAM formats more accessible and inclusive to a larger scientific demographic with off-target effects of making Python-based genomic analysis more flexible and resilient across platforms. Seqlogo was intentionally designed to decrease technical barriers of entry towards the complex concepts of sequence motifs and logo and putting the learning back into the hands of students.

Furthermore, the popularity and spread of these tools outside of their original design reveal that strategic and purposeful development of tools that diminish the user's technological burden and enrich understanding of the underlying concepts make better tools in the long run. BAMnostic and seqlogo also exhibit that there are both a need and place for graduate-level manipulatives as graduate-students are still students nonetheless and continue to need pedagogical support.

Chapter 3 Bioinformatics Interventions

An academic “intervention” is defined as “the active involvement of school officials and teachers in developing and implementing an effective plan for assisting students with academic difficulties.”¹¹⁵ Whereas manipulatives put the learning in the hands of students, interventions are systemic modifications to either the classroom, curriculum, or content that instructors use to attenuate educational barriers. As described by Lynch in 2019, an intervention should be defined by four characteristics: 1) proactive, 2) intentional, 3) formal, and 4) flexible¹¹⁶. Additionally, an intervention should also be an educational multiplier, meaning that interventions should target the skills necessary for students to further interact successfully with continued education¹¹⁷.

Students do not always progress at the same speed—a well-known and studied phenomenon which has led to the development of manipulatives (Chapter 2) and policies^{9–11} to ease the apparent differences in educational pace. The core of these approaches is to address educational inconsistencies between students to allow for accessibility to learning that is as equitable as possible. In chapter 2, we discussed the development of manipulatives specifically to aid in student understanding of complex bioinformatics concepts. The purpose of this chapter is to investigate bioinformatics curricula and suggest potential curricular redesigns as a means of academic intervention.

3.1 Threshold concepts background

In engineering, if a component fails and consequently causes the entire system to fail, it is known as a “single point of failure (SPOF).”¹¹⁸ For example, imagine a parachute; the single point of failure on early parachutes was its means of deployment—that is, the “rip cord.” If the

rip cord failed, the entire parachute system failed. Therefore, it is important to interrogate a system to identify potential points of failure and either design redundancies or develop better support structures for the potential point of failure. A similar approach can be used in education when performing curriculum evaluation during curriculum development¹¹⁹.

To illustrate the concept of curriculum failure points, visualize an international student taking two different math tests in a non-native country. The first test is comprised of only straightforward math questions, and the student passes the test. The second test is comprised of only word problems in their non-native language, and the student fails the test. Is the potential educational point of failure the expected level of numeracy or literacy? In this regard, the student outcome may not necessarily reflect the effectiveness of the instructor. Additionally, this illustration reveals that a lack of foundational understanding may also be masked by those same student outcomes. That is, if a student is successful in class, one may assume that the student fully understands the content.

A “threshold concept” (TC) is a type of curriculum failure point. “Thresholds” are portals or barriers from one area to another. A TC fundamentally alters how one perceives and understands a discipline once understood, but a TC is also a major obstacle to understanding the discipline completely⁵⁵⁻⁶³. Therefore, a TC is a concept that becomes a bottleneck to mastery of a discipline. For example, consider the concept of biological variation as it relates to the discipline of biology^{120,121}. Once fully understood, biological variation fundamentally changes how biologists perceive subjects like genetic diseases, antibiotic resistance, or virus variants. Other biology TCs like biological information, homeostasis, and evolution have been identified for the express purposes of enriching curricula with these topics to enable students more pedagogical support in understanding them¹²¹. The tradeoff of this pedagogical investment is that while an

instructor may spend less time on content along the periphery or breadth of the curriculum, the students will be given the opportunity to explore and internalize a more fundamental understanding of the discipline as a whole and accelerate peripheral comprehension through association and transfer learning⁶¹.

Because of the academic incongruity inherent to interdisciplinary science—namely bioinformatics (section 1.2)—we attempted to identify potential bioinformatics TCs with the aim of developing curricula restructuring suggestions to aid in graduate student onboarding. We explain below how a TC is defined and identified.

3.1.1 Features of threshold concepts

TCs were originally defined by Meyers and Land⁶¹ after they observed that certain concepts were believed to be fundamental for the mastery of economics and that each of these concepts has eight common characteristics^{61,122}: 1) transformative, 2) troublesome, 3) irreversible, 4) integrative, 5) bounded, 6) discursive, 7) reconstitutive, and 8) liminal. These characteristics have been used to identify and evaluate TCs across various disciplines^{55,56,58,59,63,121}.

A concept is **transformative** if—when understood—it transforms how the student perceives the discipline. That is, they undergo “both an ontological as well as a conceptual shift.”⁶¹

The **troublesome** feature is probably the most readily apparent component feature of a TC. Troublesome knowledge is defined as knowledge that is potentially “counter-intuitive, alien, tacit, ritualized, inert, [or] conceptually difficult.”⁶¹

The dissonance between student and teacher understanding of a concept exists because of the **irreversibility** of a TC. Once a TC is cognized, it is difficult to “unlearn.”⁶¹ Consider riding a

bicycle; there is the moment one does not know how to ride a bicycle and then the moment they do know how to ride a bike. Once learned (and outside of exceptional circumstances), the person can never “unlearn” how to ride a bicycle. Irreversibility also leads to the side effect called the “curse of knowledge,” when an instructor has difficulty empathizing with a student’s misunderstanding because the instructor does not remember what it was like to not know.

A concept is considered **integrative** when, upon learning, the student makes associations and connections between apparently disparate aspects of the discipline that were previously perceived as unrelated⁶¹. For example, when a student understands object-oriented programming¹²³, the component parts of a data visualization library (figure, axes, patches, etc.) become more readily understood and utilized by the student.

One of the most important, albeit abstract, features of a TC is that of its **bounded** nature. A bounded concept is a concept that is specific to a given discipline⁶¹. For example, biological variation is specific to biology and understanding it does not affect how a student perceives the field of chemistry or physics any better. As a means of foreshadowing, the bounded feature of TCs is integral to the study of bioinformatics—or any other interdisciplinary science field.

Discursive concepts are characterized by an extended or enhanced use of language relating to the discipline⁶¹. In other words, when someone studying statistics internalizes the concept of randomness¹²⁴, the word “random” takes on new meaning and fundamentally changes how the student uses that word from then on.

Reconstitutive is a somewhat meta characteristic in which, through the transformative and discursive aspects of TCs, the student undergoes a shift in subjectivity. While TC cognition is initially recognized by how one speaks and writes, it also takes place over time^{61,122}. This

gradual change ultimately reconfigures the conceptual schemas and causes both an ontological and epistemological shift¹²².

The final feature is that of **liminality**. Whereas irreversibility can lead to the “curse of knowledge,” liminality is the foundation of both “imposter syndrome” and the “Dunning-Kruger effect” because TCs are subject to three phases of progress: pre-liminal (no understanding), liminal (process of understanding), to post-liminal (understood). It is this passage between phases that can lead a student to alternate between the effects of imposter syndrome and Dunning-Kruger effect as they wrestle with understanding a TC^{122,125}.

Only with all eight features can a concept be considered a true TC. By identifying a TC, an instructor can focus on more impactful or meaningful concepts within their curriculum or begin to enrich their current curriculum with more content to support a known TC. Therefore, the identification of TCs and subsequent restructuring of curricula to support known TCs can be considered as an educational intervention.

3.1.2 Threshold concepts in graduate education

As described above, threshold concepts are the gateways to discipline-specific mastery, which is often achieved through graduate education. Therefore, identifying TCs within a given discipline serves to positively impact graduate education and research: better students, better scientists. Another unique characteristic of TCs and graduate school is that both emphasize that the rates in which students internalize and develop the content knowledge of a specific discipline are not equal. For example, two graduate students within the same department may not graduate at the same time.

TCs are similar to Piaget’s formal operational stage of development⁸⁵ defined as the point when an individual can begin to think abstractly and hypothetically. However, a conceptual

comparison between two individuals within this stage of development presents that ability as a spectrum: one person may use abstract and hypothetical thought more readily than another because of difference in needs. Likewise, the unique liminal spectrum of TCs implies that some students may fully internalize a TC and reach a post-liminal stage while other students may stay in the pre-liminal or liminal stages since only a cursory understanding of the TC may be necessary for their work or research.

An additional challenge to identifying TCs in graduate education is innovation. The landscape of curricula is constantly changing as advances in content knowledge and novelty are adopted and characterized. This innovation also extends to the differentiation of students: two students in the same department are often researching different topics, each with unique content knowledge requirements. Despite this, the bounded nature of TCs not only serve to shape a specific discipline, but also distinguish the discipline from others⁵⁹. Therefore, it is imperative that the TCs of disciplines at the graduate level be identified not only to assist in curricula development but also to inform institutional and infrastructure design.

Last, TC studies have focused primarily on more basic sciences (biology¹²⁰, chemistry⁵⁵, statistics¹²⁴, and computer science¹²³) with some recent work in medically-related fields^{60,62,63}, whereas burgeoning or interdisciplinary fields are woefully understudied^{55,59}—specifically at the graduate-level. The source of this disparity in coverage is likely due to a number of factors like research vs pedagogical interests (section 1) and smaller class sizes. Educational research for a specific discipline is likely to be reviewed and published within the same discipline. This suggests that educational research for a given discipline is often assessed based on the discipline's conventional research norms (e.g., large population studies and quantitative data), making it difficult to publish impactful research in relevant interdisciplinary journals.

3.1.3 Interdisciplinary science threshold concepts

The most prominent impediment to identifying TCs within an interdisciplinary field is the bounded feature of TCs; a TC is discipline specific. “Interdisciplinary” is taken to mean “any form of dialogue or interaction between two or more disciplines.¹²⁶” The claim is that the outcome or use of the amalgamation of sub-disciplines within an interdisciplinary field are fundamentally different than that of the sub-disciplines by themselves. That is, “interdisciplinary outcomes do not represent the sum of the constituent disciplines⁵⁹.” Bioinformatics is considered an interdisciplinary field made up of computer science, mathematics/statistics, and biology^{44,46}. Therefore, all content knowledge and conceptual understanding of bioinformatics is either a 1) collection of disparate concepts from its constituent sub-fields or 2) collections of concepts of its constituent sub-fields that—when combined—are unique to bioinformatics.

Therefore, we asked whether it is possible to identify a bioinformatics-specific TC or whether any perceived bioinformatics TC is potentially a combination of TCs from computer science, biology, and mathematics/statistics and not bioinformatics specific.

3.2 Identification of threshold concepts within bioinformatics

For context of our study, all data were collected from self-selected students and faculty affiliated with the Department of Computational Medicine and Bioinformatics (DCMB) at the University of Michigan in late 2021. Data from this study may include responses from undergraduate students as part of the Advanced Master’s Degree Program (AMDP). The bioinformatics curricula include components of computer science, mathematics/statistics, and biology and is considered an interdisciplinary science with the common focus of creating novel informatic and computational methods, tools, and algorithms for basic biomedical, translational, and clinical research. The approximate student base at the time of this study was 180 students

(88 PhD, 76 MS, and 16 AMDP). According to the Integrated Postsecondary Education Data System (IPEDS), there were only 86 institutions in the United States that awarded bioinformatics degrees in 2020¹²⁷. The purpose of this study was to conduct surveys and focus groups of both DCMB-affiliated students and faculty to ascertain whether bioinformatics-specific TCs exist, and, if so, potentially identify some initial TCs to better inform bioinformatics curricula development and design.

3.2.1 Bioinformatics threshold concept study design

Since our work was educational research and our survey data was anonymized and aggregated, we applied and were approved for institutional review board exemption (HUM#00185545). The irreversible nature of TCs (e.g., “curse of knowledge”) implied that the work should be student-centered. In other words, initial surveys and data collection should be driven by responses and feedback from proximal learners, an approach unlike most contemporary research as those studies developed prompts for students based on faculty input first^{55,59,60,62,63}.

The study was divided into five phases (**Figure 11**). To generate focus group and survey prompts for proximal learners, we performed a review of bioinformatic pedagogical literature and syllabi from within DCMB (Phase 0). The data from the review were collected based on emphasis of subject (e.g., length of time the topic is covered) or depth of coverage for a given topic (e.g., the number of times a topic was addressed). These data were coded, pruned, and categorized using the constant comparative approach^{128,129} to develop discussion and survey prompts ubiquitous enough that proximal learners would likely recognize the core concept and provide feedback appropriate to the modality in which it was presented.

Student focus groups and DCMB student-wide surveys (Phase 1) would provide the bulk of the data to begin generating a list of problematic curricula content, as perceived by the students. These data would then be presented to a faculty focus group (Phase 2) to revise and refine the list of potential TCs. If necessary, we would recruit faculty from comparable bioinformatic programs (Phase 3) to further revise and refine the list and continue to iterate as necessary (Phase 4+).

Identification of threshold concepts (TC) within bioinformatics (HUM#00185545)

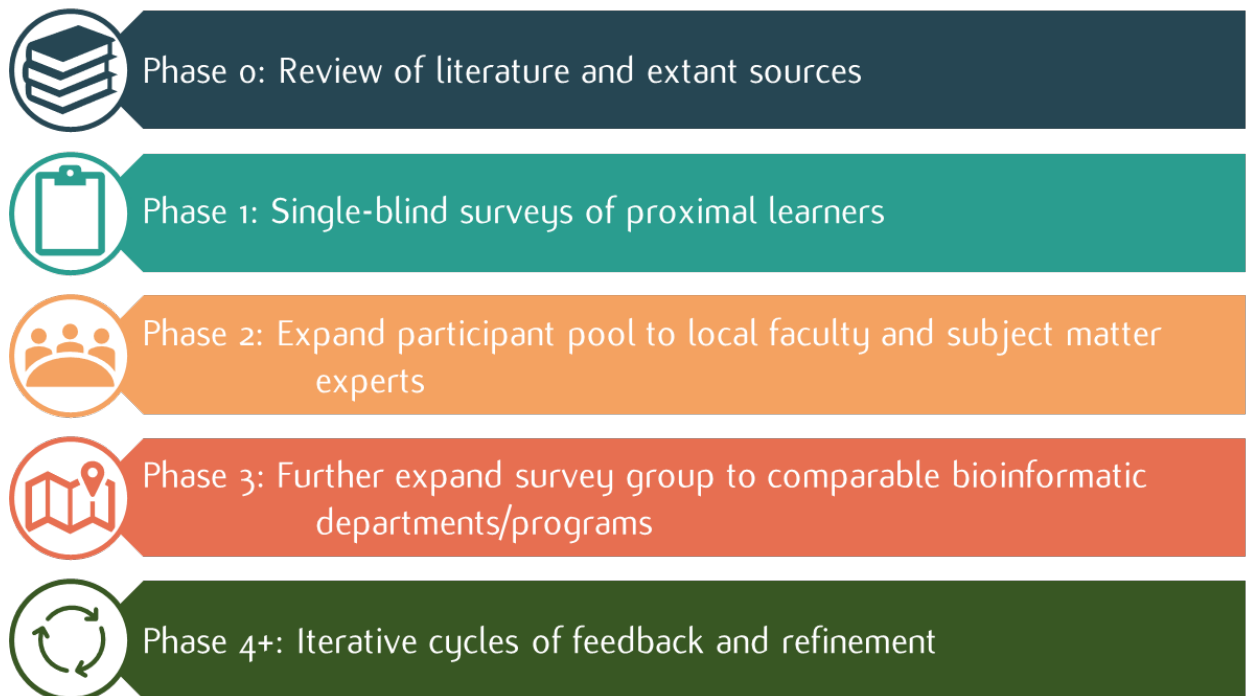


Figure 11 Bioinformatics threshold concepts study design

The original design of this study comprised 5 phases. Phase 0 was an initial review of literature and extant sources (e.g., syllabi). Phase 1 was single-blind surveys to proximal learners. Phase 2 recruited local faculty and subject matter experts to refine and revise original data obtained from Phase 1. Phase 3 was meant to recruit comparable bioinformatics programs and further refine and revise the list of potential TCs. Last, Phase 4+ indicated iterative cycles of feedback and refinement.

3.2.2 Bioinformatics student focus group

Using the discussion prompts from Phase 0, a student focus group was formed. In 2021, an invitation to a “Bioinformatics academic innovation focus group” was sent to the student body of DCMB. Due to the COVID-19 pandemic, this meeting was managed virtually via Zoom. The invitation required only that participants be affiliated with DCMB. Approximately 20 participants took part in a two-hour conference call that contained multiple small group breakout sessions and helped identify an initial list of “problematic” concepts from the bioinformatic curricula.

The focus group, however, was not without limitations. The first major limitation was that the inherent nature of remote conferencing, while potentially confidential, precluded anonymity. This limitation was mediated by not recording the session and taking notes based only on general discussion from the group. This led to the second limitation in that it was difficult to maintain record keeping due to paraphrasing by small group representatives or overall discussions over specific topics. Last, and most important, was the limitation of self-selection of the participants. Since the modality of the focus group was discussion-driven, there was a selection bias for fluent English speakers. While international students may have significant interest in impacting or providing feedback regarding academic innovation, remote conferencing can impact international student participation¹³⁰.

After weighting the discussion prompts based on the initial student focus groups, we decided to mitigate selection bias and increase response rates by developing a text-based DCMB student-wide survey.

3.2.3 *Bioinformatics student-centered survey*

The DCMB student-wide survey was developed using the Qualtrics platform¹³¹. We leveraged the platform's ability to anonymize respondent submissions while fully managing deployment and collection. The survey was comprised of 49 questions (**Appendix A**). Question types include 5-point Likert scale arrays (e.g., strongly disagree, disagree, neutral, agree, and strongly agree), item ranking, best option selection, and free text entry. Participants could choose to answer any, all, or none of the questions as they saw fit, including demographics. The survey was divided into five components: demographics, biological concepts appropriate to bioinformatics, mathematical and statistical concepts appropriate to bioinformatics, computer science concepts appropriate to bioinformatics, and general open-ended questions. Finally, the survey was pretested via Qualtrics platform tools but no additional or intentionally designed quality improvement was implemented.

The invitation to participate was delivered through a DCMB-affiliated student listserv. Participants had approximately one month to respond (between 9/27/2021-10/24/2021) and were reminded to participate every week until the survey was closed. No incentives were offered to participants. We received 70 responses to the survey ($\approx 40\%$ response rate based on total DCMB-affiliated students). Responses were cross-referenced by their self-reported length of time in the DCMB program (**Table 3**). Of the 70 responses, 48 indicated that DCMB was their primary department. Additionally, students were stratified by their self-reported primary field of research (**Figure 5**). Non-response bias was not assessed, and no validity framework¹³² was implemented.

<i>Length of time in program</i>	<i>%</i>	<i>Count</i>
< 1 year	34.38%	22
1-2 years	28.13%	18
2-3 years	14.06%	9
3-4 years	12.50%	8
4+ years	10.94%	7

Table 3 Self-reported time in program from DCMB student survey

3.2.4 Bioinformatics faculty focus groups and survey

Several small DCMB faculty focus groups were conducted to refine the feedback from DCMB student focus groups. The faculty focus groups took place via remote conferencing in late 2021. The participating faculty were given context to the educational goals of the research and then prompted to discuss the student feedback and categorize and/or generalize a given concept or topic. All faculty participants were self-selected from an open invitation to all DCMB-affiliated faculty via faculty listserv.

After completion of the DCMB student-wide survey, a DCMB-affiliated faculty survey was developed to provide final feedback on the student-provided list of conceptually difficult topics within the bioinformatics curricula. The survey was created and pretested using the Qualtrics platform¹³¹. The survey was 57 questions comprised of 5-point Likert scale (e.g., strongly disagree, disagree, neutral, agree, and strongly agree), best option selection, and free text entry (**Appendix B**). Participants could choose to answer any, all, or none of the questions as they saw fit, including demographics. The survey was divided into seven sections: demographics, sequential data analysis, mathematics/statistics, data management, uses of extant knowledge (e.g., references and calibration curves), data scaling, and open-ended free text. Aside from the demographics and open-ended sections, each section was modeled the same: a concept was introduced; faculty were asked whether the concept any of the following applied: fundamental to bioinformatics, troublesome to understand, unique to bioinformatics, and transformative; and faculty were asked to provide corroborating observations or feedback regarding the concept.

The invitation to participate was delivered through a DCMB-affiliated faculty listserv. Participants had approximately two weeks to respond (between 1/21/2022-2/7/2021) and

reminded to participate every week until the survey was closed. No incentives were offered to participants. We received 19 responses (\approx 53% response rate out of 36 primary and joint faculty). Faculty were asked to self-select their association with DCMB as either “primary faculty” (52%), “joint faculty” (21%), or “research faculty” (5%). Non-response bias was not assessed, and no validity framework¹³² was implemented.

3.2.5 Limitations of data collection and analysis

Because bioinformatics is a broad interdisciplinary field, there is no standard sequence of curricula. Education at the graduate level also includes specialization training outside the scope of this study. Additionally, since this study was developed from a student-centered position, troublesome concepts were provided by the students and refined by faculty. Therefore, unlike other TC studies^{55,120,123}, it was impractical to develop low-stakes assessments for students with the scope of the complete bioinformatics knowledgebase because a specific student’s specialization may preclude them from adequately addressing topics outside of their specialty. That is, a molecular dynamics student may not need to know everything required of a genomics specialization. This suggests that developing assessments prior to identifying the fundamental concepts of bioinformatics would only cause undue pressure on both the study and participants. Likewise, the lack of performance-based assessment also precludes a pre-post comparative analysis as both the anonymous nature of the results and time between assessments make this type of analysis intractable and, ultimately, outside the scope of the study. Last, student feedback could likely be influenced by courses most recently undertaken either because the recent study made fundamental knowledge gaps salient or because it increased perceived competence in a fundamental skill.

3.3 Results

This section presents the analysis of the findings, DCMB student focus groups and surveys, and DCBM-affiliated faculty focus groups and surveys. The purpose of this analysis was to identify whether TCs exist within the bioinformatics discipline and—if they do—how they can be used to inform further bioinformatics curricula design and restructuring. The results from the student (n=70) and faculty (n=19) surveys provide most of the qualitative data while the focus groups allow for more prosaic descriptions of difficult bioinformatic concepts. It is important to note that at no point was a topic triaged or prescribed by the researchers. Therefore, all topics addressed by students or faculty were dependent on the topic's prevalence throughout preceding phases of the study. Thereby, a topic may have been presented to a further phase regardless of whether it met the criteria of a TC solely to act as a prompt for further discussion or become subject to refinement.

Given the features of TCs, our approach to identifying potential bioinformatics TCs was divided into two components: 1) isolate concepts that DCMB students had difficulty processing early in the program but took for granted towards the end of the program and 2) leverage faculty insight in determining whether a potential TC were troublesome to student, transformative to understanding bioinformatics, unique to bioinformatics, and fundamental to bioinformatics. To accomplish the first component, we cross-referenced student responses with length of time within the program divided into 5 groups: <1 year, 1-2 years, 2-3 years, 3-4 years, and 4+ years. A newer student should respond less confidently to TCs than more senior students. Additionally, since TCs are notably troublesome, a specific concept may not actually be grasped by students until much later; therefore, concepts uniformly difficult across all demographics were also

flagged as potential TCs. Concepts consistently difficult to newer students and concepts difficult at all levels were presented to the faculty focus groups and surveys.

3.3.1 Student focus group

Review of extant literature and curricula sources identified several recurring bioinformatic topics that required regular remediation. These topics were used as prompts for DCMB student focus groups. The topics emphasized or emphatically discussed the most by and between DCMB students can be found in **Table 4**. Many of the topics identified by students were linked to computation-based concepts. This list of troublesome bioinformatics topics was subject to an initial review by a DCMB faculty focus group. At this point in the study, the faculty were not provided the TC context of the study and were just asked to refine the language or potentially recategorize the topics. This refined list can be found in **Table 5**.

Troublesome bioinformatics topics identified from student focus groups

How to identify statistical distributions that apply to given biological data
Various characteristics of high-throughput data and analysis (e.g., error-correction and biases)
How to “start” a problem
Lazy vs eager loading
How to properly debug and troubleshoot
Indexing strategies
Sequence alignment algorithms
Tool development and deployment

Table 4 Troublesome bioinformatics topics identified by student focus groups

<i>Faculty-refined list of troublesome topics provided by DCMB students</i>	
<i>Core tasks</i>	
	How to start problem solving
	Central dogma
<i>Application of concepts</i>	
	How to understand and when to apply statistical properties to identify statistical distributions
<i>Domain-specific</i>	
	Multiple sequence alignment
	Sequence similarity searching
	Functional motif searching
	Structure prediction
	Bioinformatics literature review
	Sequence assembly
	Sequence alignment and their respective algorithms
	Reference genomes
<i>Programming in bioinformatics</i>	
	Indexing strategies
	Tool development and deployment
	Debugging strategies
	Code review
	Lazy vs eager loading
	Data structures
	HPC & High-throughput data

Table 5 Faculty-refined list of troublesome topics provided by DCMB students

3.3.2 Student survey

The DCMB student survey collected 70 responses. This sample size, while comparatively small, is representative of the DCMB student body ($\approx 40\%$ response rate) and consistent with other contemporary research in comparable fields^{55,59,123}. As mentioned, the purpose of this survey was to elicit responses regarding the refined troublesome bioinformatics topics in **Table 5**. The students were asked to respond based on level of agreement with qualitative questions (**Appendix A**) regarding the topics. These data were then cross-referenced based on length of time within the program to identify topics that junior students responded to less confidently than more senior students. Since there are fewer senior students than junior students and, therefore, fewer respondents of each demographic, we normalized responses based on number within a given demographic.

The three most prominent examples of troublesome bioinformatics concepts—as identified by student responses and Welch’s 2-tailed t-test (**Table 6**)—were 1) determining statistical properties of bioinformatic data ($p < 0.05$, **Figure 12**), 2) the concept and application of references to extant knowledge ($p < 0.01$, **Figure 13**), and 3) debugging strategies ($p < 0.005$, **Figure 14**).

Computing sequence similarities of bioinformatic data ($p = 0.7$, **Figure 15**), lazy vs eager loading ($p = 0.69$, not pictured), and leveraging indexing strategies ($p = 0.9$, not pictured) are unlike the concepts above since they presented with somewhat uniform confidence among all demographics. Ideally, senior students should always be more confident with a troublesome topic relative to junior students. However, both student groups present with similar middling confidences across these troublesome topics. This suggests that these topics could be more difficult to internalize in general, thus more senior students still have not passed that given

threshold of understanding. This implication may potentially indicate the presence of a higher order TCs as well.

Student free-text submissions were divided into three sections: biology, mathematics/statistics, and computer science. In each section, the students were provided an opportunity to describe any relevant bioinformatics concept or topic not covered for the given section. Within the biology free-text entries, the most enriched topic was related to sequencing technologies (e.g., Single cell RNA-seq, ATAC-seq, etc.). There were no significantly enriched topics within the mathematics/statistics free-text entries and the provided entries spanned from “reading and understanding mathematic formulae” to Bayesian inference and parameter estimation. The computer science free-text entries were enriched for topics of parallelization and machine/deep learning.

After processing student survey response data, the following list of student-reported curricula obstacles was generated and used to prompt refinement and revision by the DCMB faculty:

- Sequential data analysis
- Statistical distribution(s) identification and application
- Data ingest, exploration, and management
 - Identifying potential edge cases
 - The principles and approaches to data cleansing
 - Data wrangling
 - Data exploration
- Data scaling
- References to extant knowledge

<i>Question</i>	<i>Junior Mean</i>	<i>Junior SD</i>	<i>Senior Mean</i>	<i>Senior SD</i>	<i>Mean Diff</i>	<i>P value</i>
<i>I understand how to perform bioinformatic literature review in my field</i>	3.27	1.25	4.17	0.69	0.89	0.05
<i>I understand the concept and application of references of extant knowledge (e.g., reference genomes, reference sets, and standard curves)</i>	3.45	0.89	4.33	0.47	0.88	0.01
<i>When necessary, I know how to interrogate bioinformatic data to determine their statistical properties</i>	3.00	0.79	4.00	0.82	1.00	0.04
<i>When necessary, I know how to identify the statistical distribution(s) that apply to a given set of bioinformatics data</i>	3.31	0.98	3.67	0.75	0.35	0.42
<i>I understand how to perform a proper code review</i>	3.06	1.09	3.60	0.49	0.54	0.17
<i>I can confidently exercise debugging strategies</i>	3.50	1.17	4.80	0.40	1.30	<0.005
<i>I know how to package and deploy a program in a scalable and maintainable way</i>	2.63	1.22	3.40	1.20	0.78	0.29
<i>I can confidently use version control</i>	2.88	1.22	4.60	0.49	1.73	< 0.001
<i>I can confidently navigate and manipulate the command line interface</i>	4.06	1.09	4.40	0.49	0.34	0.38
<i>I can quickly deconstruct a programmatic task into its constitutive components</i>	3.63	1.11	4.20	0.75	0.58	0.25
<i>I understand the difference between "lazy" and "eager" loading with respect to data analysis</i>	2.00	1.13	1.80	0.75	-0.20	0.69
<i>I understand when and why certain data structures (e.g., dictionaries and dataframes) are used for data analysis</i>	3.79	1.01	4.60	0.49	0.81	0.05
<i>I know how to access and leverage multiple indexing strategies (the representation of an item's position within a sequence) across multiple data structures</i>	3.50	1.12	3.40	1.36	-0.10	0.90
<i>I comprehend the concepts, implementation, and application of dynamic programming algorithms</i>	3.07	1.22	3.60	1.02	0.53	0.41
<i>I know how to scale prototype algorithms to high-performance and high-throughput computing</i>	2.21	1.26	3.20	1.47	0.99	0.27
<i>I comprehend functional motif searching and identification</i>	2.93	1.28	3.40	1.36	0.47	0.56
<i>I can confidently compute sequence similarities of two (2) or more sequences</i>	3.50	1.12	3.20	1.83	-0.30	0.77

Table 6 Welch's 2-tailed t-test between junior and senior student responses

Bolded p values indicate a lack of differentiation in confidences of a given topic between demographics.

Level of agreement with "When necessary, I know how to interrogate bioinformatic data to determine their statistical properties"
(n=56)

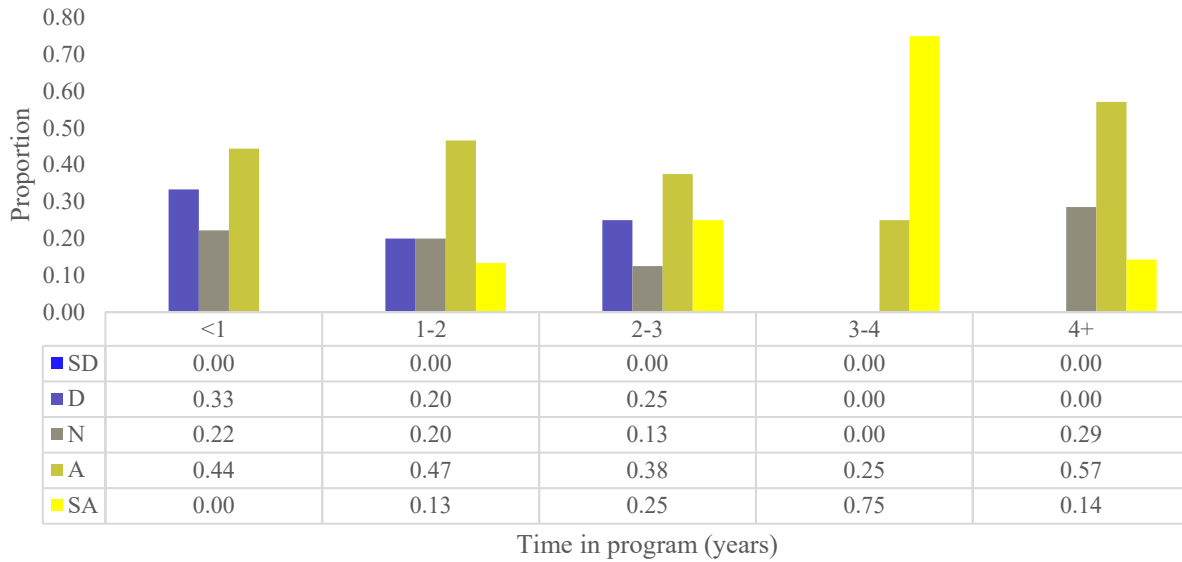


Figure 12 DCMB student feedback on statistical properties of bioinformatic data (n=56)

Junior students (less than two years) generally responded less confidently on how to determine statistical properties of bioinformatic data than their senior counterparts.

Level of agreement with "I understand the concept and application of references to extant knowledge"

(n=63)

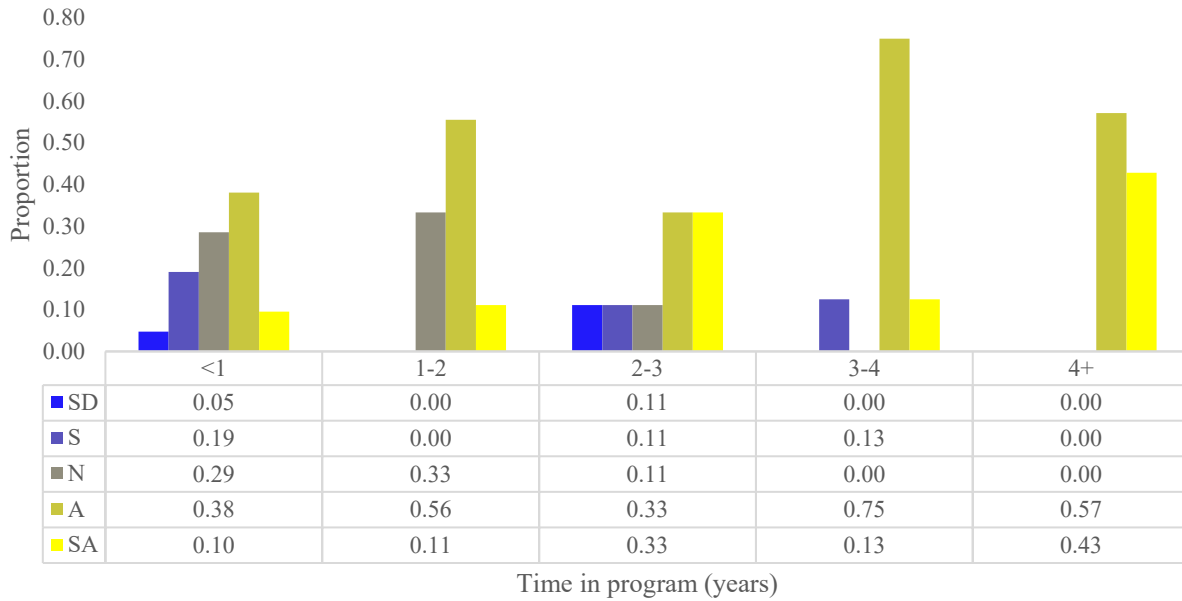


Figure 13 DCMB student feedback on references to extant knowledge (n=63)

Junior students (less than two years) generally responded less confidently on the application of references to extant knowledge (e.g., reference genomes and calibration curves) than their senior counterparts.

Level of agreement with "I can confidently exercise debugging strategies"
(n=54)

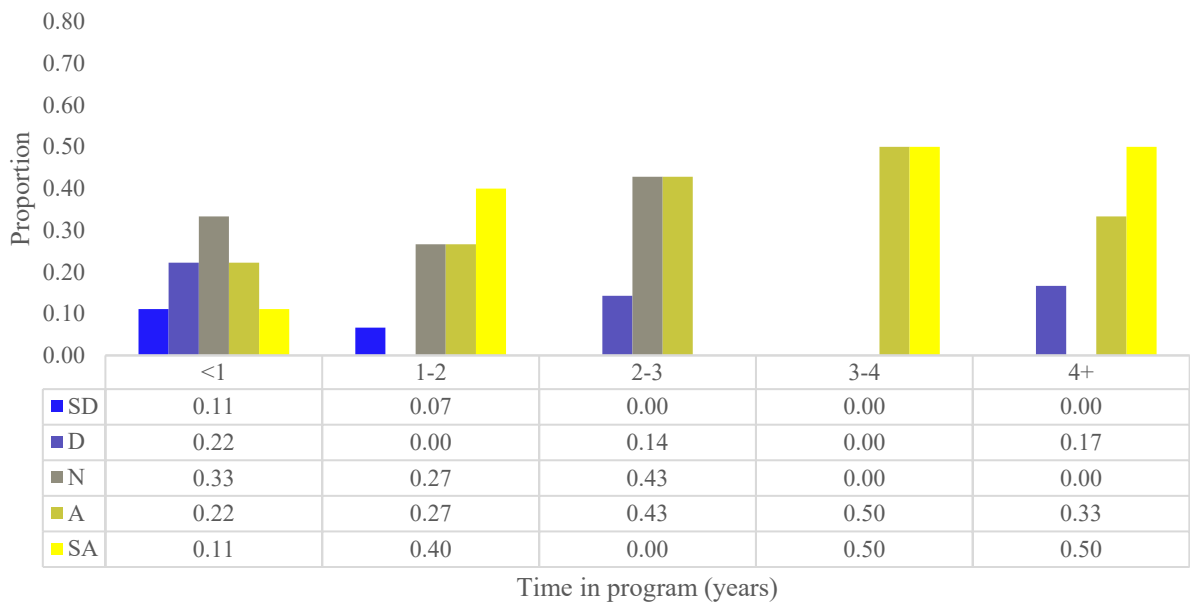


Figure 14 DCMB student feedback on debugging strategies (n=54)

Junior students (less than two years) generally responded less confidently on how to exercise appropriate and sustainable debugging strategies for bioinformatic tools than their senior counterparts.

Level of agreement with "I can confidently compute sequence similarities of two (2) or more sequences"

(n=52)

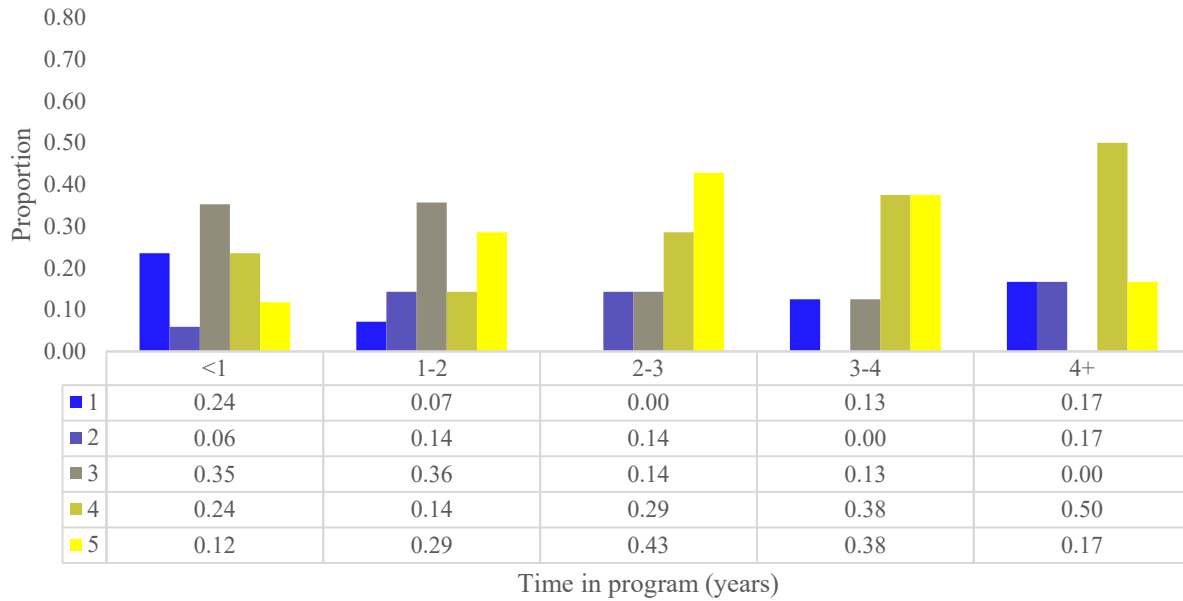


Figure 15 DCMB student feedback on sequence similarities of bioinformatic data (n=52)

Junior students (less than two years) generally responded less confidently on how to determine statistical properties of bioinformatic data than their senior counterparts.

3.3.3 Faculty feedback

The DCMB faculty survey collected 19 responses. The purpose of this survey was to ask for insight into identifying potential learning obstacles of bioinformatics—in general—within our bioinformatics program. The faculty were asked to respond based on level of agreement with qualitative questions (**Appendix B**) regarding the student-identified bioinformatics topics at the end of section 3.3.2: sequential data analysis; statistical distribution(s) identification and application; data ingest, exploration, and management; references to extant knowledge; and data scaling.

The survey then divided each of the topics in the student-identified list into individual sections. Each section asked faculty to reflect on whether the given topic was 1) transformative to how the field of bioinformatics is understood, 2) troublesome to understand, 3) specific to bioinformatics, and 4) fundamental to understanding bioinformatics as a field. Depending on the insight of the faculty, potential TCs were identified if they met all four criteria. We represented the data using a diverging stacked bar chart to visualize the overall sentiment of the topics (**Figure 16**). Faculty generally agreed that each topic presented in the survey was transformative and fundamental to bioinformatics. With the minor exception of the topic of sequential data analysis (where faculty were split), each topic was also generally agreed upon as being troublesome to understand. Finally, faculty generally disagreed that each topic presented was specific to bioinformatics.

DCMB-affiliated faculty survey of potential threshold concepts
(n=13)

■ Strongly disagree ■ Disagree ■ Neutral ■ Agree ■ Strongly agree

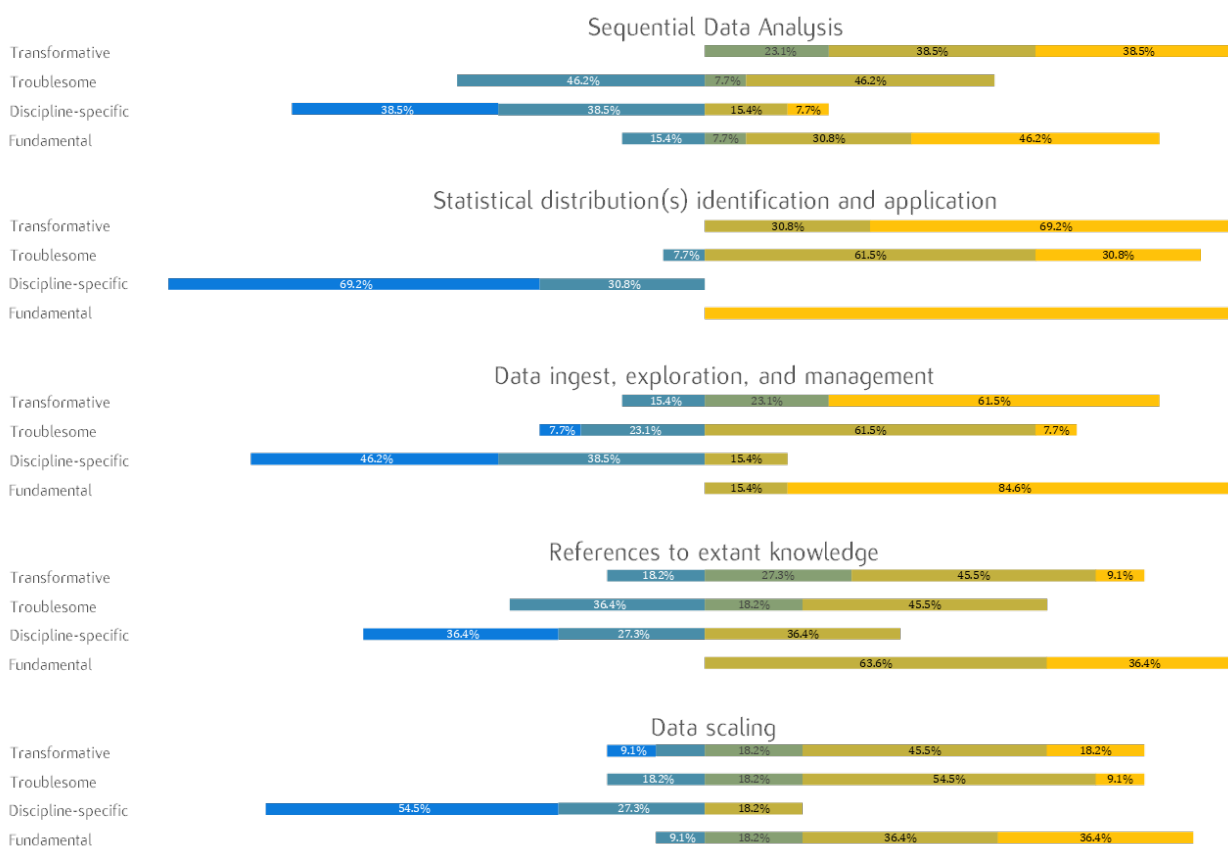


Figure 16 DCMB-affiliated faculty survey of potential threshold concepts (n=13).

Faculty generally agreed that every topic presented to them was transformative, troublesome, and fundamental. However, faculty generally disagreed that any topic was bioinformatics-specific.

3.4 Discussion

To inform the development of any academic interventions in bioinformatics, we conducted a study to identify potential TCs that may be present in bioinformatics curricula. It has been suggested that there is a possibility that the interdisciplinarity of a field produces unique outcomes that are unlike the sum of its constituent disciplines^{55,59}. Given that unique academic environment of graduate education within the highly interdisciplinary field of bioinformatics, we first had to determine whether TCs exist within the field. To accomplish our study, we performed several student and faculty focus groups and surveys to gain insight into known or perceived conceptual obstacles within the DCMB bioinformatics curricula.

Early in the study, we postulated that student feedback would strongly be enriched for computational concepts. This assumption is based on the understanding that the common core between all students within DCMB may not be the biological content they study—as students likely have varying specializations—but rather the computational approaches utilized within the field. Of the five most prominent bioinformatics concepts generated by student feedback, three of them related significantly to computational approaches. Student feedback strongly suggested a general malaise or perceived dissonance between their ability and departmental expectations. This suggests that DCMB students recognize an apparent deficit in their understanding of computational concepts and require additional remediation, recitation, or differentiation not currently offered within the department or readily available outside of the department.

Ultimately, faculty feedback regarding the student-generated list of conceptual stumbling blocks suggested that none of the topics met the criteria of a threshold concept. That is, while most of the faculty agreed that the concepts were transformative, troublesome, and fundamental

to bioinformatics, the faculty also generally believed that none of the topics was specific to bioinformatics. This finding agrees with results from other research in interdisciplinary fields⁵⁹.

Threshold concepts are difficult to identify. Faculty often have forgotten what it was like to not know (curse of knowledge), and students do not know what they do not know (black swan). Hypothetically, a concept may appear as a threshold concept at face value due to a spurious observation caused by a more fundamental concept that is yet to be fully understood. This study may not have identified or disproven the existence of threshold concepts in bioinformatics. It did, however, expose some potential avenues of academic intervention and curricula retuning within bioinformatics.

The first avenue for academic intervention and development is that early computational onboarding of incoming students is taught primarily through the lens of the Python and R programming languages (**Figure 3**), but several special topic courses offered within the DCMB utilize more disparate languages (e.g., MATLAB and Julia). Therefore, it may be advantageous for the department to encourage or enforce instruction to continue utilizing the preexisting scaffolding developed by the required courses students have already taken. This standardized programming language approach could potentially offload the technological barriers to learning experienced by students when exploring more advanced bioinformatics concepts and specializations within the department.

The second (and potentially the most pernicious and difficult to address) avenue for academic intervention development is that of academic siloing. At the undergraduate level, education is designed to span multiple disciplines by using prerequisites. A computer science student may take biology and political science outside of the computer science department. However, graduate education within biomedical/interdisciplinary sciences often does not follow

that same path. If a bioinformatics student wants to take additional computer science classes, many of the fundamental computer science courses are outside their reach. Therefore, departments often take stop-gap measures to address the need by creating courses within the department. Consequently, these courses can often be more broad or topical than necessary, ultimately lacking in deeper fundamentals than temporal restrictions allow. Both academic and logistical issues could be addressed by intentionally de-siloing curriculum and creating avenues of collaboration between these departments. For example, computer science departments could develop courses for bioinformatics students and vice versa with bioinformatics developing courses for computer science students. Only by diversifying thought does a multidisciplinary institution become more agile and resilient.

Last, when considering the results of this study, it was decided that Phases 3 and 4+ (**Figure 11**) would not be necessary. The original design of the study assumed that a list of locally defined potential bioinformatics TCs would be defined and subject to a more global refinement. Since there is no locally defined list of potential TCs, this study was closed.

Chapter 4 Pedagogical Professional Development

4.1 Introduction

One of the major obstacles to onboarding new graduate students is that many programs lack the pedagogical scaffolding to effectively advance borderline or disenfranchised students quickly. The most straightforward path to overcoming this obstacle is to improve the quality of the teaching. As reviewed in chapter 1, faculty at research-focused institutions are recruited to interdisciplinary fields based primarily on the impact of their research and less so on formal educational training. Consequently, while these fields are taught by some of the best scientists in their fields, the scientists are not taught how to teach effectively. The lack of pedagogical emphasis is further perpetuated and compounded by perceived cultural norms and selective pressures towards research. We further interrogated DCMB faculty survey data described in section 3.2.4 to potentially substantiate whether bioinformatics faculty lacked pedagogical emphasis and professional development.

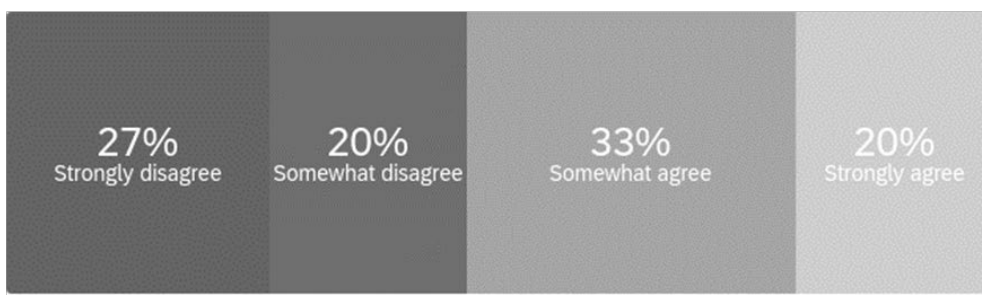
4.1.1 Faculty Survey

In early 2022, an open invitation to participate in a DCMB-affiliated faculty survey regarding the study of bioinformatics threshold concepts, we received 19 responses. The survey is detailed further in section 3.2.4 and the questions can be found in **Appendix B**. Faculty were not incentivized to respond, and participants could choose to answer any, all, or none of the questions; all submissions were anonymous. For the purposes of this chapter, we explored only the initial demographic data collected from the survey. Of the responses, over 50% of the

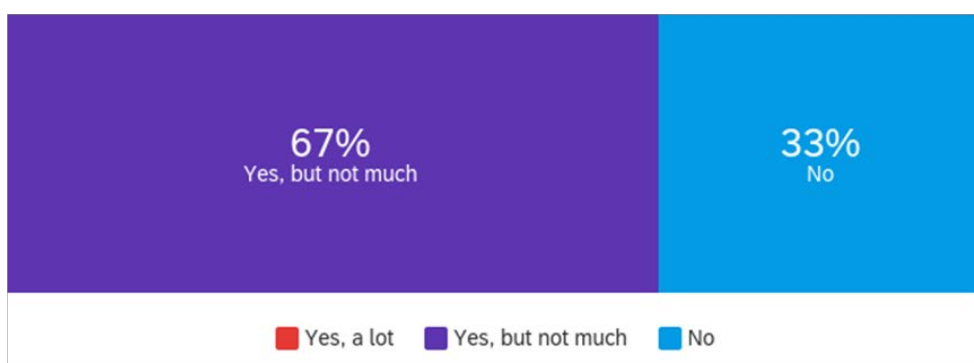
respondents reported having a didactic, or classroom, teaching effort of 15% or more. 33% of the respondents reported having no formal pedagogical training while 0% reported having a significant amount. Formal pedagogical training was defined as workshops, courses, or seminars specializing in pedagogical training. Last, 60% of the faculty respondents reported having 20+ student contact hours per year. (**Figure 16**). While these data are not powerful enough to determine statistical significance, they connote a lack of pedagogical emphasis and professional development among bioinformatics faculty of DCMB—albeit qualitatively.

With the assertion of the presence of pedagogical biases among the faculty, we postulated that the current pedagogical framework of DCMB serves as a proxy for other biomedical science departments and programs and that those faculty match that of the cultural evolutionary theory model detailed in section 1.2.3. The conceptual model describes the transmission of pedagogical biases within academia (**Figure 4**)⁷⁶. Evaluation of that model suggested that the three most impactful areas to introduce and incentivize pedagogical training were: 1) PhD and postdoctoral students who are likely to become research-focused faculty, 2) faculty that train graduate students in research-focused institutions, and 3) the channels between teaching- and research-focused institutions⁷⁶. To ameliorate this pedagogical bias, we launched a pilot training program directed toward current graduate students within the biomedical sciences of the University of Michigan Medical School (UMMS).

A I consider 15% or more of my effort to be dedicated to didactic teaching?



B Do you have any formal pedagogical training?



C Estimate your number of contact hours with students per year

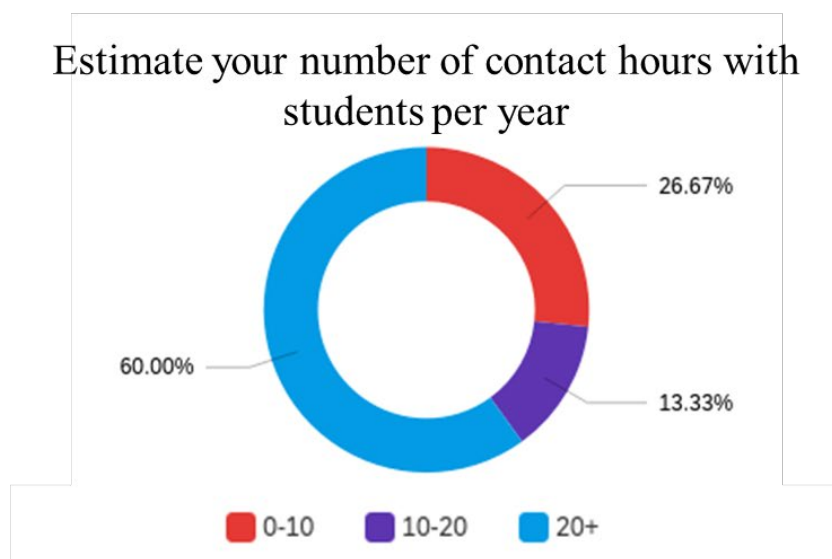


Figure 17 DCMB-affiliated faculty survey demographic data (n=19)

A. The stratification of faculty responses based on percent effort dedicated towards didactic teaching. **B.** Self-identified level of formal pedagogical training of faculty respondents. Formal pedagogical training was defined as programs such as "seminars on lesson planning, inclusive teaching, or preparing future faculty workshops." **C.** Representation of faculty respondents' estimation of student contact hours per year.

4.2 Pedagogy of Interdisciplinary Science Education

The effectiveness of an educator is parameterized by three types of knowledge: subject matter knowledge, pedagogical knowledge, and pedagogical content knowledge (PCK). Pedagogical content knowledge (PCK) is the unique understanding and approaches a teacher uses to support a student's learning of a subject matter^{37,38,67}. Videlicet, a science teacher must know science, how to teach, and how to teach science. Forasmuch as graduate students are arrayed with subject matter knowledge, pedagogical knowledge and PCK are the most impactful targets of pedagogical professional development. We applied for and were awarded a fellowship for the pilot cohort of UMMS' Research | Innovation | Scholarship | Education (RISE) program. Through innovation coaching and exercising our community of practice made available by RISE, we developed and operationalized the Pedagogy of Interdisciplinary Science Education (POISE) training program (<https://poise.med.umich.edu/>) to address the gap in pedagogical methodologies, approaches, and knowledge contextualized through interdisciplinary science content.

The two-fold purpose of POISE is to teach the future teachers within biomedical interdisciplinary science how to effectively teach and to answer the question "Does pedagogical training for the instructors improve the quality of the students?" POISE partnered with the Center for Research on Learning and Teaching (CRLT), Center for Academic Innovation (CAI), and the School of Education to curate a curriculum that covers formal pedagogical training all governed within the context of how it applies to both graduate-level education and the interdisciplinary sciences, in other words, how to effectively teach graduate-level, interdisciplinary scientists. To pursue POISE's purpose, the pilot cohort was comprised of students who had some prior graduate-level teaching experience in their fields

4.2.1 Theory of Change

Early in POISE's development, we outlined a theory of change (ToC). A ToC is graphical representation of why a theory will work and by what means it can be empirically measured^{133,134}. ToCs are a product of program evaluation techniques for theory-driven approaches¹³³⁻¹³⁶. A program's ToC is co-created through collaboration between the program team and interested stakeholders and subject to multiple rounds of iteration before a program is fully developed¹³⁴. A traditional ToC is comprised of five components: 1) impact, 2) outcomes, 3) outputs, 4) activities, and 5) inputs. This listing also suggests another aspect of ToCs: they follow the backward design framework²⁹. Each component will be explored through the context of POISE's ToC (**Figure 18**).

The *impact* of a ToC is the expected long-term systemic change correlated with a program's output¹³⁷. Impact is often the first thing to be identified during program development. For example, POISE's impact was defined as the following (red box in bottom right of **Figure 18**): "Formal pedagogical training for instructors (of any role) in the biomedical sciences is encouraged, protected, and incentivized at the university, school, and departmental levels." This is a long-term approach to shift professional biases within the biomedical sciences towards a community of practice that supports and incentivizes pedagogical professional development in accordance with cultural evolution theory; students who expect and respect well-trained educators at the graduate-level will themselves become well-trained educators to future students. This premise serves as a positive feedback loop that could potentially shift academic cultural norms and values.

Any *outcome* of a ToC is designed in tandem with outcome *indicators*. An outcome (short- or mid-term) represents the intended and unintended consequences of the program's

operation¹³⁷. An indicator is the means in which the outcome is assessed, either qualitatively or quantitatively. POISE's outcomes are found in the bottom center portion of **Figure 18** represented by red (outcomes) and green (indicators). A short-term outcome of POISE was that trainees would experience "elevated professional self-confidence," and the indicator was that trainees would present with improved teaching evaluation results following their participation in POISE.

A ToC *output* is defined as the immediate products of the program's operation and are necessary for achieving outcomes¹³⁷. An output can serve as a measure of the program's cadence and progress. POISE's outputs are defined (bottom left of **Figure 18**) as number of applicants (proxy for interest), cohort retention rate, comparative analytics (cohort progress), committee feedback (governance and operationalization), results of low-stakes assessments, and self-/peer-evaluation.

The *activities* of a ToC are the direct means in which the program enacts change for their outcomes to be realized¹³⁷. POISE activities (blue rectangle in top left of **Figure 18**) are seminars, workshops, application (mock-lectures), networking, and practicums. In education, a "practicum" often takes the form of embedding a trainee with an experienced educator and having them critically observe the educator on topics like classroom management and student participation.

Last, the *inputs* of a ToC are resources devoted to or promised to the program to ensure the outcomes are achieved¹³⁷. POISE's resources (green iconography in top right of **Figure 18**) are training from program staff, funding from RISE, and support from POISE partners. Additionally, POISE describes the *issue* that it is attempting to address in conjunction with the *impact* it ultimately is attempting to realized. The issue ties to the general theme of this

dissertation: “impactful student understanding and training, especially in graduate-level biomedical science, is often impeded due to a lack of formal pedagogical training of the instructors.”

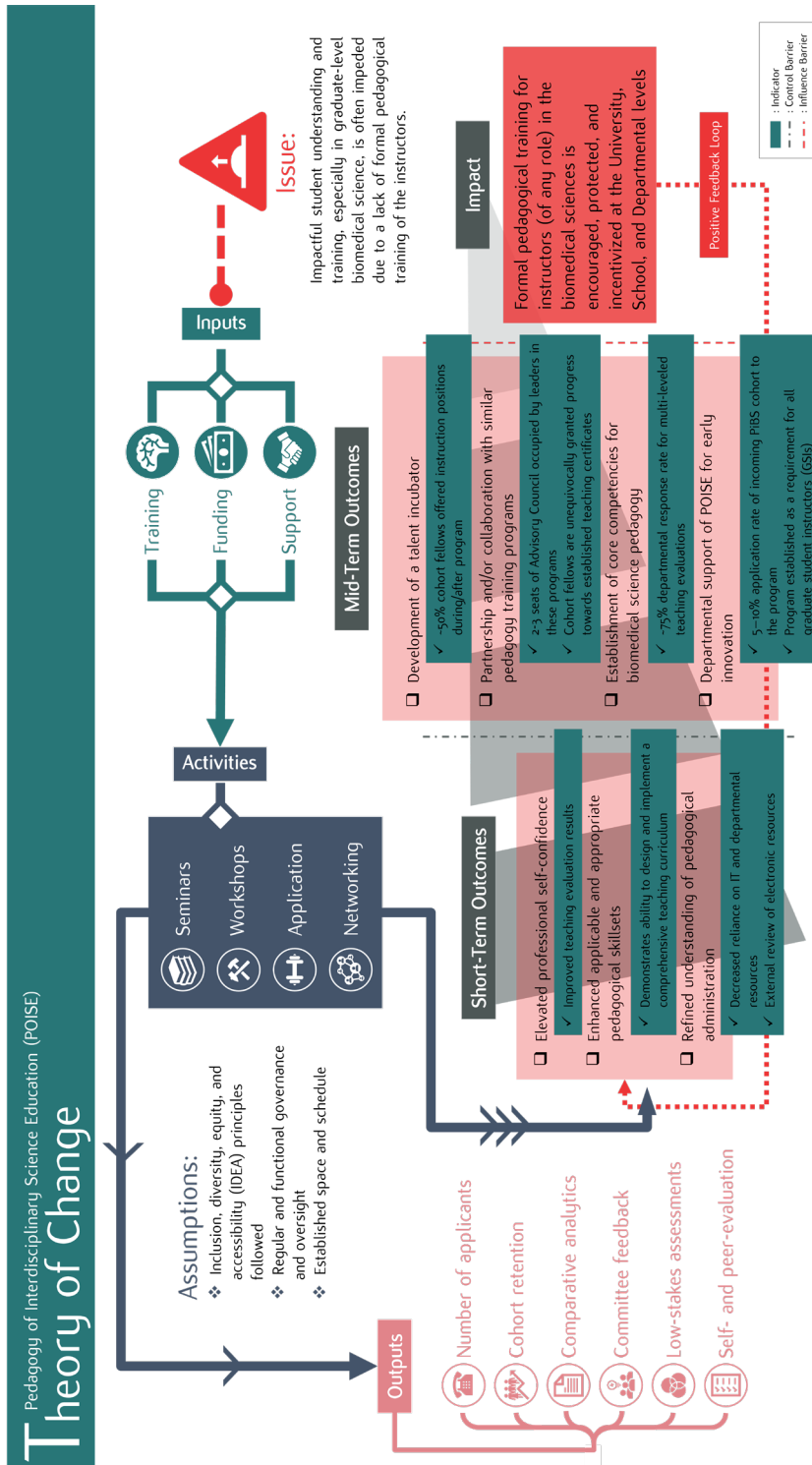


Figure 18 POISE Theory of Change

The Theory of Change is a representation of the intended impact, outcomes, outputs, activities, and inputs of POISE to address the observed issue.

4.2.2 Development and design

Innovation is both the creation of something new as well as the application of preexisting solutions to new problems in different ways. POISE sought to innovate education in the interdisciplinary sciences through the introduction of classical pedagogical methodologies present in education curricula that state-licensed educators must follow for their degrees—albeit in a much more abbreviated and curated form in POISE—all the while being contextualized to the interdisciplinary sciences. The reason the latter concept is important is that developing a lesson plan for a class on chromosomal architecture to genetics students is completely different from developing a lesson plan for a class on programmatic detection and categorization of chromosomal architecture states to genomics students. One is just the biology. The other is the biology plus the computer science and statistics. Further attenuating retention and comprehension is that students new to these interdisciplinary fields tend to come from institutions that do not offer specialized interdisciplinary undergraduate programs. This means that the “consumer to producer” paradigm of graduate-higher education is diminished since the students in these fields are constantly cycling back to a consumer role more during production phases. To guide our curriculum, POISE was designed around six core components (**Figure 19**): learning theories and education frameworks, classroom administration, advanced instructional methodologies, professional development, practicums, and application.

The means by which the program would be contextualized to the interdisciplinary sciences would be dictated by the pilot cohort, possible because in conjunction with the content, the trainees were cumulatively creating a capstone project that would demonstrate the practical application of their content within their field of research. The trainees would be required to present and periodically peer-review each other’s projects throughout the program and speak to

the projects' application and efficacy. For the capstone project, each trainee was tasked to identify a course they had taken within their field of research. This course had to be a course that was difficult not because of the content, but because of how the content was delivered. Once the course was identified, the trainee would utilize the theories, methodologies, and educational techniques they were introduced to during training to redesign their identified course. This process required trainees to create and manage classrooms on learning management systems (e.g., Canvas), develop syllabi, write contribution guidelines and teaching philosophy statements, design a lesson plan with full scope and sequence, and generate a full lecture/class with an accompanying form of student assessment. The final student artifacts took the form of individually identified and designed courses that culminated in a peer-reviewed mock-lecture.

POISE Design Components

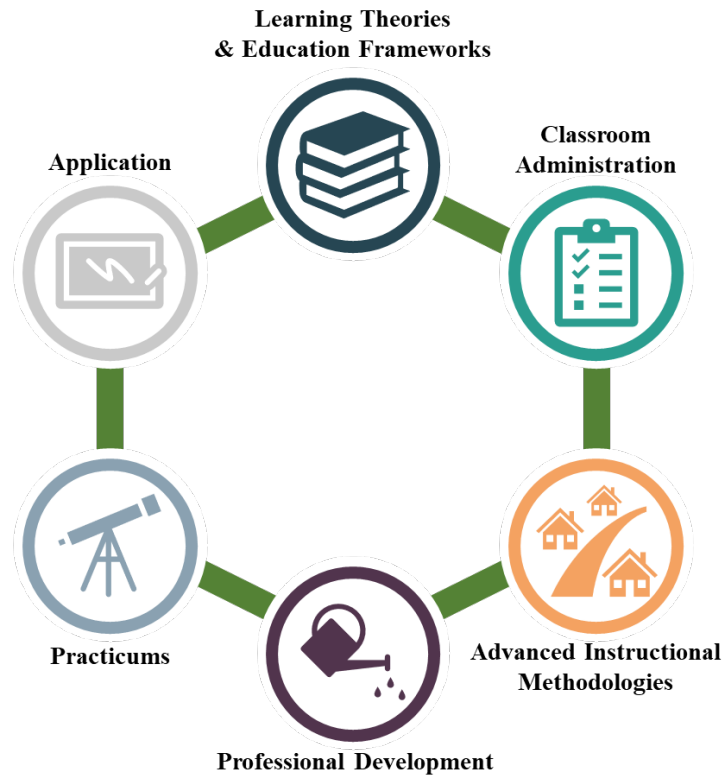


Figure 19 POISE design components

POISE is divided into six core components: 1) learning theories and educational frameworks, 2) classroom administration, 3) advanced instructional methodologies, 4) professional development, 5) practicum, and 6) application.

4.2.3 POISE Advisory council

To ensure that POISE had guidance and oversight, an integral part of POISE's development was the identification and recruitment of invested stakeholders to act as the POISE advisory council (PAC). To that end, we exercised the RISE community of practice to recruit our PAC members. The PAC was composed of members who could provide insight into each of the components of POISE. Don Peurach, PhD from the University of Michigan's School of Education and Tazin Daniels, PhD from the University of Michigan's Center for Research on Learning and Teaching (CRLT) joined to speak on the educational and teaching components of POISE. Ryan Mills, PhD from DCMB and Scott Barolo, PhD from the University of Michigan's Program in the Biomedical Sciences (PIBS) provided support through their interdisciplinary science acumen and expertise. Rachel Niemer, PhD from the University of Michigan's Center for Academic Innovation (CAI) and Rajesh Mangrulkar, MD; Paula Ross, PhD; and Nikki Zaidi, PhD from RISE provided guidance and backing regarding the operationalization, development, and innovation of POISE.

The PAC met every one to two months during the development of POISE, applications process, and throughout the training cycle of the pilot cohort. The PAC provided advice and guidance on everything from application questions to the order of the curriculum and was instrumental to the success of POISE's pilot.

4.2.4 Curriculum

In accordance with the six core design components of POISE (section 4.2.2), we divided our curriculum into three phases: 1) theoretical knowledge, 2) applied knowledge, and 3) practical knowledge (**Figure 20**).



Figure 20 POISE curriculum

The POISE curriculum is divided into three phases. The first phase (Theoretical Knowledge) covers theories on how learning occurs, and teaching frameworks designed to support learning. Phase 2 (Applied Knowledge) covers the application of pedagogical knowledge to support classroom administration such as lesson planning and rubric design. Phase 3 (Practical Knowledge) introduces advanced educational methodologies such as active learning and inclusive teaching.

Phase 1 (Theoretical Knowledge) of the POISE curriculum primarily covered learning theories and educational methodologies. The purpose of this phase was to instruct the cohort on how learning happens. To develop an understanding of the ways in which students both learn and develop, we introduced Piaget's Theory of Cognitive Development⁸⁵, Sweller's Cognitive Load Theory⁶⁵, and Vygotsky's Social Development Theory and Zone of Proximal Development¹³⁸. The topics of Bloom's Taxonomy^{30,139} and Threshold Concepts^{55-63,120-125} (Chapter 3) were dedicated to explicating guided and advanced learning processes. Last, we clarified core competencies^{5,140,141} and how they can be used in competency-based education.

Phase 2 (Applied Knowledge) covered administrative techniques utilized by more classically trained educators. The purpose of this phase was to develop how teaching should be developed. Chief among these concepts was the introduction of the Backwards Design Framework²⁹ which explains the process of creating a program, curricula, or even lesson plans by first identifying the desired achievable result, establishing appropriate and authentic means of assessing whether the result were met, and only then creating the content that supports the accomplishment of result. In that regard, the Backwards Design framework is very similar to the ToC (4.2.1). With this context, more general approaches like creating teaching objectives and lesson planning were delineated. To assist in the evaluation step of backwards design, specific focus was given to the development of authentic student assessment and evaluation, rubric design, and how to define and weight appropriate classroom participation and attendance¹⁴⁰⁻¹⁴⁶.

Phase 3 (Practical Knowledge) presented advanced instructional methodologies (e.g., think-pair-share and active learning) and contemporarily informed concepts (e.g., inclusive teaching and trauma-informed teaching). The purpose of this phase was to train how to use the training within the context of current political, sociological, and emotional landscape. Advanced

technological techniques like virtual classroom management and the role of technology in the classroom were included, especially within the context and importance of the COVID-19 pandemic taking place during POISE's operation.

4.2.5 Pilot cohort

POISE originally budgeted for a pilot cohort of five trainees. This budget was dictated by costs of manipulative development and distribution, space rentals, refreshments, and travel costs for trainees to attend guided practicums with local educators. However, an early stay-at-home order at the start of the COVID-19 pandemic was issued, and the budget would therefore welcome a much larger cohort given the virtual requirements of the training. The call for applications was issued to all graduate students within UMMS in October of 2020.

We received 17 responses. Of those 17 responses, 11 were accepted into the POISE training program (7 females, 4 males; 9 within their 4th year, 1 in their 1st, and 1 in their 2nd) that represented seven different programs or departments within UMMS. The selection criteria were as follows: 1) that participants had either completed at least one term as graduate student instructor (GSI) or were currently a GSI, 2) that participants had clear intention to pursue a teaching role following graduation (e.g., professor, trainer, etc.), and 3) that participants had support from their research advisors. Over the course of POISE's training (December 2020-May 2021), we had a 100% cohort retention rate with all trainees successfully completing every activity and requirement of the curriculum.

4.2.6 Program evaluation

Program evaluation was determined by a retrospective pre-post analysis similar to Kirkpatrick's program evaluation model^{135,136,147-153}. Upon completion of the training, we

collected 11 exit surveys (100% response rate). The survey was comprised of 30 questions (**Appendix C**). Question types include 5-point Likert scale arrays (e.g., strongly disagree to strongly agree), multiple choice, item ranking, and free text. All answers were anonymously collected through the Qualtrics platform¹³¹. Non-response bias was not assessed, and no validity framework¹³² was implemented. No incentives were offered to participants.

The primary portion of the pre-then-post data was collected from question 4 of the survey (5-point Likert scale array). Trainees were prompted to reflect on the topics and modules covered over the course of their POISE training. The prompt requested they also reflect on those same topics from a perspective prior to POISE. Additionally, all free text answer to open-ended questions were coded for sentiment and constant comparative analysis^{128,129} using the Qualtrics toolkit.

4.3 Results

After collecting trainee responses, we performed a 2-tailed t-test of the retrospective pre-post data ($df = 10$, **Figure 21**). All responses indicated a statistically significant difference between the two groups. The least significant ($p = 0.013$) was observed regarding their agreement with how confidently they could manage a virtual classroom. The most significant ($p < 0.001$) were observed for their understanding of learning theories and how to apply learning theories to develop curriculum.

Additionally, participants were asked to rank each of the topics covered throughout POISE so that we could identify impactful modules to expand upon later. We generated a diverging stacked bar chart to visualize POISE topic rankings (**Figure 22**).

Retrospective Pre-then-Post Evaluation

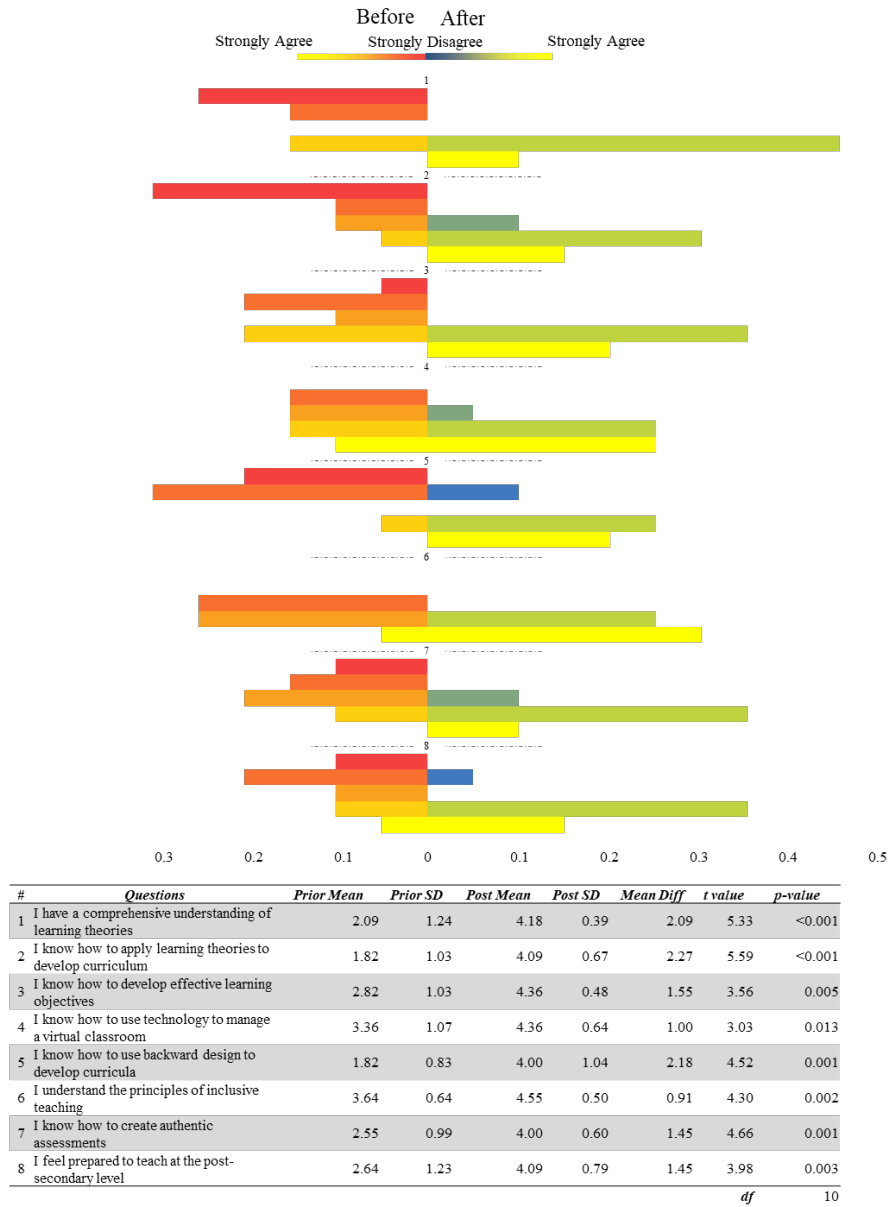


Figure 21 Retrospective pre-post evaluation

Using trainee exit surveys, a 2-tailed t-test was applied to trainee responses regarding agreement and confidence levels prior to and directly following POISE training. The top portion of the figure plots agreement and confidence levels (yellow representing strong agreement) for both prior to (left) and following (right) POISE training. The bottom table contains the 2-tailed t-test data. All questions presented as significantly statistically different ($p < .05$), with the largest observed differences relating to understanding and application of learning theories ($p < 0.001$) and the smallest observed difference relating to technology in the classroom ($p = 0.013$).

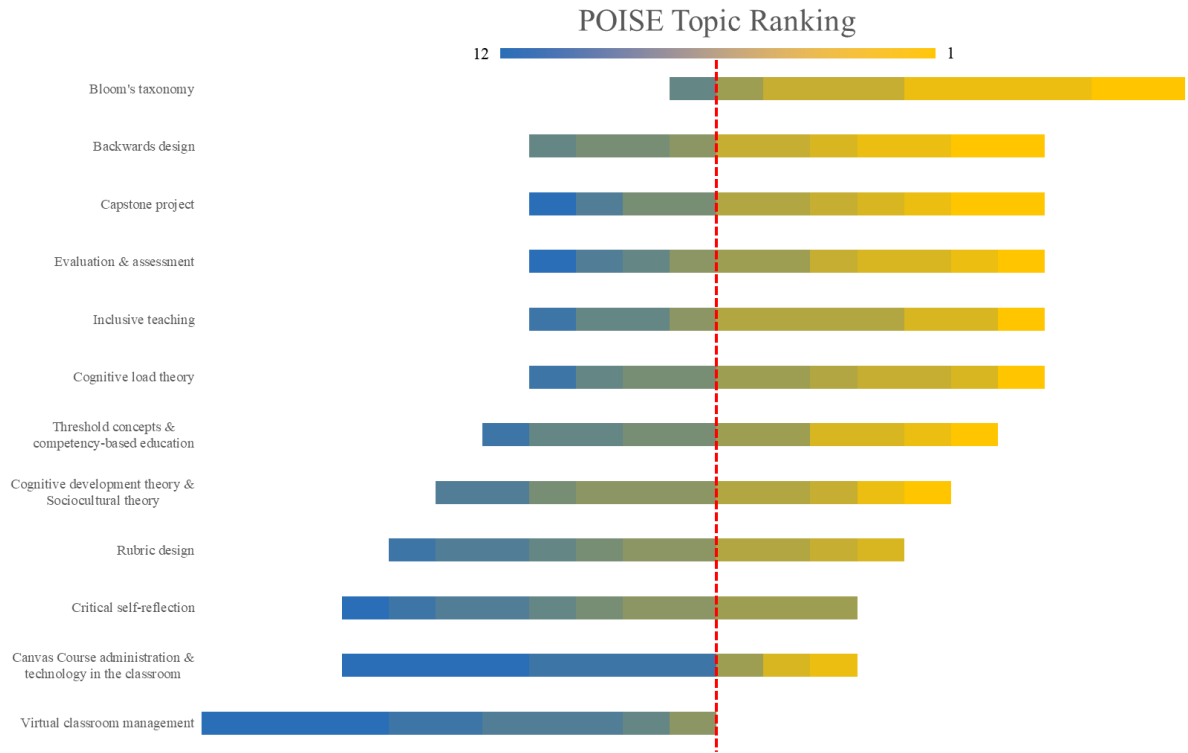


Figure 22 POISE topic ranking

Diverging stacked bar chart of each of the 12 formal topics covered within the POISE curriculum. The color of a discrete box within a given topic represents its rank (1 = yellow, blue = 12) where rank 1 represents the highest. The size of the discrete box represents the proportion of responses. These charge does not suggest a topic's necessity, but rather its perceived weight relative to the other topics. The highest rank topic was Bloom's taxonomy while the lowest was virtual classroom management.

4.4 Discussion

The cynosure of this chapter was teaching how to teach effectively. We believe that the most formative development for instructors is achieved by refining their pedagogical skill sets (PCK). Teaching, however, is not just the content covered, but also how content is taught. This hypothesis reflects the apparent gap present within graduate level interdisciplinary science education. The faculty are thoroughly subject matter experts. Like students, however, they may experience certain academic deficiencies at this level in that while students may lack the content knowledge, faculty may lack the pedagogical knowledge to effectively convey the content to the students. Given selective pressures within modern academia (Section 1.2.3 **Figure 4**), we investigated the effect of pedagogical professional development by targeting the future educators and proximal learners. To essay this investigation, we developed a pilot program (POISE) with the purpose of teaching the future teachers within biomedical interdisciplinary science how to teach effectively.

Through the POISE pilot cohort, we demonstrated a significant change in confidence and understanding across a wide range of pedagogical skill sets (**Figure 21**). The largest differences in confidence were observed through topics related to learning theories and education frameworks. Free text responses indicated a revelatory response when topics like Bloom's Taxonomy³⁰ were covered because it informed the participants on how to create actionable and impactful lessons that supported the evaluation and assessment techniques necessary for backwards design lesson planning²⁹. Another informative observation on technology was evident in both confidence and topic ranking (**Figure 22**); the participants were less influenced by topics related to technology in the classroom or virtual classroom management. As mentioned earlier, the entirety of the POISE pilot cohort was operationalized during the COVID-19 pandemic.

Therefore, many of the trainees were inundated with remote teaching and research. We hypothesize that this overexposure had a desensitizing effect on trainees, making the topics carry less weight. The most important result was that the trainees had a significant change in pedagogical confidence. While completely anecdotal, nothing influences the impact of an educator more than their own self-confidence.

While this research attempted to answer whether pedagogical training for the instructors improves the quality of the students, the work was primarily to innovate in the space of graduate education to address the apparent gap in pedagogical training in interdisciplinary science. There were many limitations to the program that should be mentioned. First was the scale's limit. As a pilot program with a cohort of 11 trainees, both statistical power and curriculum development required more input. Put simply, this was intended as proof of principle and would require additional iterations to fully understand the impact of the program or properly shape a fully informed curriculum. Additionally, given the timeline and constraints of the COVID-19 pandemic, we are unable to directly answer whether or not our training has had downstream effects on students. The last limitation of POISE is ownership. This limitation stems from the question of whom POISE belongs to and, therefore, who funds it. Again, owing to the COVID-19 pandemic, institutional resources are understandably focused on sustainment and support. Only as the academic infrastructure adapts and stabilizes following COVID-19 will the continuity and succession of POISE be addressed.

It would be disingenuous, however, to suggest that programs like POISE do not already exist. One of the questions we received regularly from stakeholders was “how is POISE different from other programs?” The most comparable resource to POISE would be that of the Center for Research on Learning and Teaching (CRLT) at the University of Michigan¹⁵⁴. This program, and

many like it, serve as monoliths that include a multitude of services from pedagogical training to curriculum evaluation. We fully acknowledge CRLT as an excellent standard to compare against for resources and training in teaching and education. A majority of the workshops, however, handle high-level educational concepts (e.g., active learning) while eschewing more fundamental ones (e.g., how to evaluate student performance) and often present in a general manner that better serves courses in the humanities through the lens of undergraduate education. CRLT, as it is, is not agile enough to readily address the needs of esoteric and diverse niches present within graduate level education. Therefore, while CRLT is serving its current audience well, it is not serving all potential audience members. This is the key difference between CRLT and POISE. POISE focuses on science education at the graduate-level. As such, POISE supplements the existing CRLT curriculum by addressing an underserved niche in academia.

POISE has also been likened to future faculty on-boarding workshops (e.g., preparing future faculty or PFF¹⁵⁵) and the Institutional Research and Academic Career Development Award (IRACDA)^{156,157} programs like the University of North Carolina at Chapel Hill's Seeding Postdoctoral Innovators in Research and Education (SPIRE)¹⁵⁸. Like POISE, PFF workshops and seminars target graduate students to provide them with teaching-related professional development, whereas IRACDA programs target postdoctoral candidates to provide them the training for both biomedical research as well as teaching. The key differences between POISE and PFF is that of time commitment, the scope and sequence of the content, and a capstone project that produces an actionable and contextualized product developed by the trainees. POISE intentionally requires participants to attend multiple sessions over an extended period of time with significant participation. Because of this requirement, POISE is enabled to explore and introduce a broader scope and sequence of content without compromising depth. The key

difference between POISE and IRACDA programs is that of accessibility. Positions within IRACDA programs are competitively awarded and funded through grants from the National Institute of Health (NIH). Furthermore, IRACDA programs are dichotomous in their nature. They are meant to support the scientists as well as teachers by enforcing an equally divided training. Finally, and probably the most deliberate difference between these other programs and POISE is that POISE contextualizes the training to the interdisciplinary sciences prior to career placement.

Chapter 5 Conclusion

I want to take a moment to reiterate and emphasize the first words of this dissertation: in no way is this dissertation meant to be a scathing indictment on the education that I have received at the University of Michigan or the Department of Computational Medicine and Bioinformatics. Quite the opposite. This dissertation was made possible only through the cooperation, support, and guidance of the faculty. Any critique or criticism is desired to be purely constructive and to ask, “how can we be better?” so that we can further support our claim as “Leaders and best.”

In chapter 1 of this dissertation, we introduced three issues we believe are affecting bioinformatics education. These three issues are the accessibility of content (section 1.2.1), effective curriculum design (section 1.2.2), and pedagogical content knowledge (1.2.3). To address these issues, we used three overarching questions to guide our work: “how do bioinformatics students learn?” (Chapter 2), “what are bioinformatics students learning?” (Chapter 3), and “who is teaching our students?” (Chapter 4). Consequently, we will conclude this work by probing into the future directions of our projects.

5.1 Bioinformatics manipulative development

In Chapter 2, we asked, “how do bioinformatics students learn?” and defined “accessibility of content” as an issue of unintentional curriculum design that resulted in constant code switching of our students. To address this issue, we developed two pedagogical manipulatives: one to aid student access to commonly used genomic data formats (BAMnostic) and another to reduce technological barriers to exploring the concepts of sequence motif

identification and visualization (seqlogo). These two programs served as a proof of concept that intentional design for the express purpose of easing access to content for the students supports student understanding of the concepts. Additionally, the same intentionality leads to a more user-friendly program in general.

These two programs were built to meet specific observed use cases, so the future of this line of inquiry is dependent on a more general approach to bioinformatics program evaluation. Additional use cases for manipulative development, however, may be identified through rigorous program evaluation. One such potential use case is that of graph-based concepts. A graph is made of nodes and edges. A node represents some data, and the edges are the relationships that node has with others. This concept is used to explore many bioinformatics concepts like de Bruijn graphs¹⁵⁹, Hidden Markov Models¹⁶⁰, phylogenetics, etc. All these concepts can be easily explored using classic object-oriented programming approaches, and therein lies the issue. Some students do not come from formal computer science backgrounds. Therefore, developing a manipulative program that serves as an abstraction layer for the students would prevent students from expending significant cognitive load on understanding the underlying concept of object-oriented programming and allow them to focus on the current content instead.

The most important aspect of our bioinformatics manipulative development was to show the usability, marketability, and impact of tools designed specifically for education within the interdisciplinary sciences. This shift in paradigm is markedly divergent from cultural norms within both biomedical research and interdisciplinary science in general, namely, a positive, lasting effect on a given field is possible without traditional research and publication modalities. While that impact can be observed through non-traditional metrics (e.g., number of downloads), the more poignant impact is on our students. By enabling our students with a more

straightforward approach to understanding, science—as a whole—is significantly and broadly impacted.

5.2 Threshold Concepts

In Chapter 3, we were guided by the question “what are bioinformatics students learning?” to interrogate the issue of ineffective curriculum design. To investigate this issue, we used a student-centered, bottom-up approach to identify difficult conceptual content bioinformatics students were facing. A major finding of this research was not that bioinformatics does not have any threshold concepts—that is yet to be proven. This study illustrated the difficulty in identifying TCs within interdisciplinary science. That is, if a concept like version control were identified as a commonly misunderstood concept among bioinformatics students does that make it a bioinformatics TC? No, because concepts like version control and statistical analysis are not unique to bioinformatics. Consequently, this characteristic of students with TC deficiencies from constituent disciplines may be a feature of the system.

By elucidating the existence of these TC deficiencies—or “orphan concepts”—a conceptual framework could be developed to outline an ontology to define interdisciplinary science (**Figure 23**). In other words, by describing the required patterns of TCs that compose an ideal scientist in an interdisciplinary science, an interdisciplinary science defines its own bounds, thereby questions like “what is the difference between biostatistics and bioinformatics?” could be answered objectively. Additionally, by creating a map of these desired TCs, orphan concepts could be intentionally addressed as a means of academic interventions (**Figure 24**). For example, incoming bioinformatics students can take a battery of low-stakes assessments. The results of those assessments could identify orphan concepts that each student is lacking, and their training

would be tailored to meet those deficiencies. This procedure would optimize graduate student learning while reducing the general pedagogical load of the faculty.

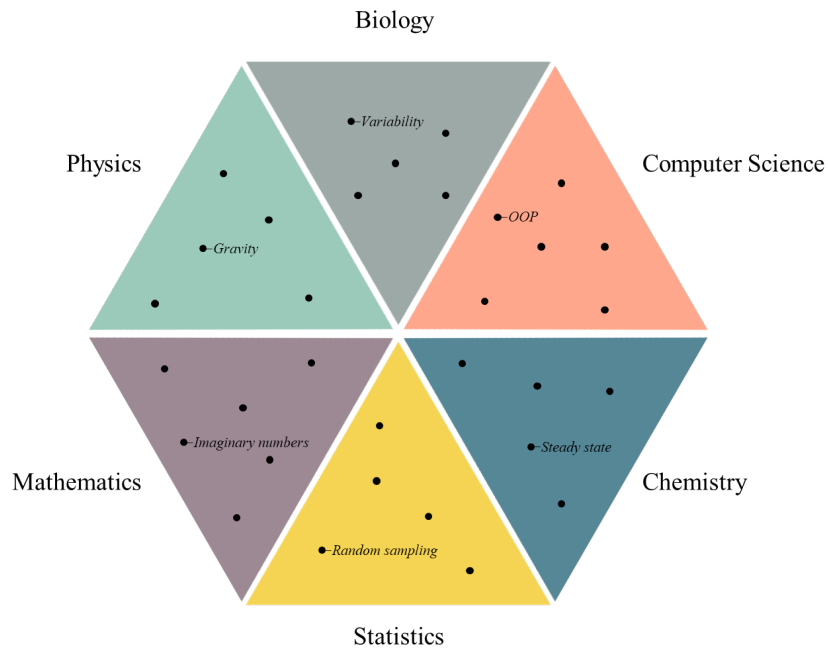
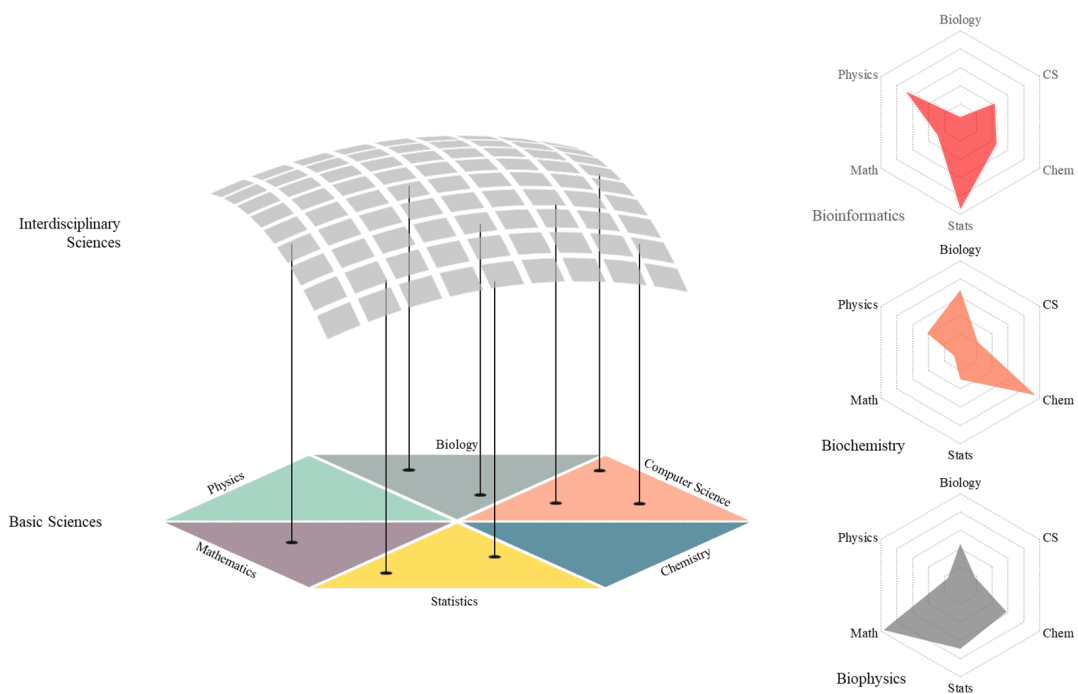
A**B**

Figure 23 Threshold concept ontology

A. Representation of some known TCs and to which discipline they belong. A feature of TCs is they are discipline specific. **B.** Suggests a representation of how interdisciplinary fields can be defined by the pattern of the TCs required from their constituent disciplines.

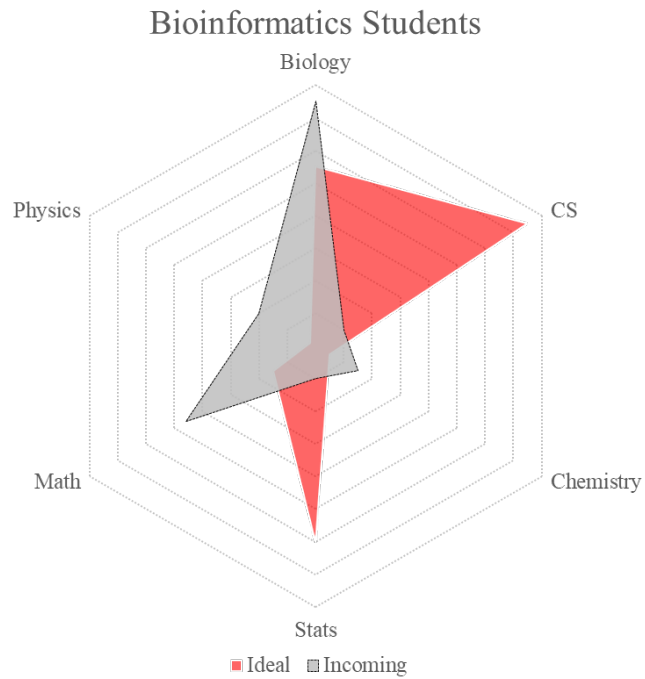
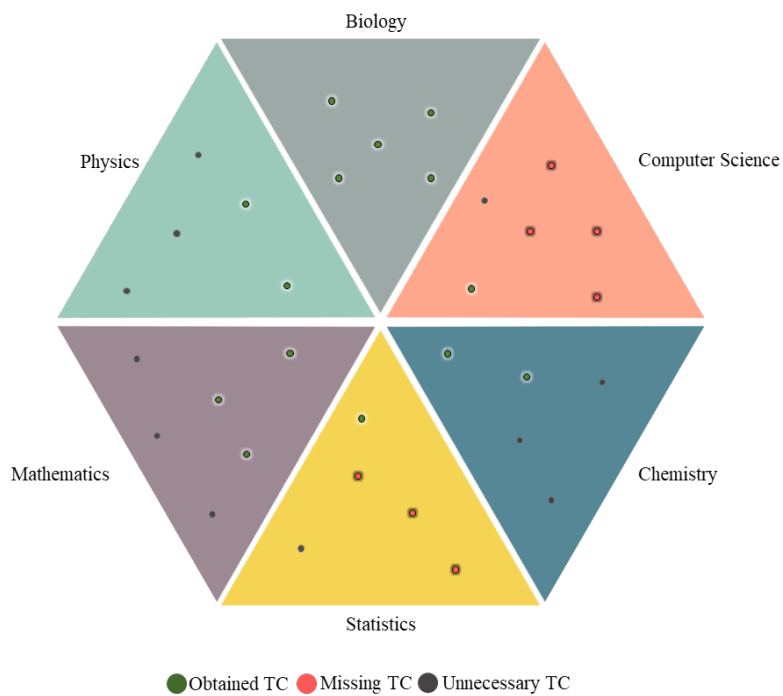
A**B**

Figure 24 “Orphan” Concepts

A. Students may join a field of interdisciplinary science with both known and unknown deficiencies. **B.** By using the ontology explained in **Figure 23**, a department can perform robust entrance assessments to determine which TCs a student may be missing and tailor student progression and training.

5.3 POISE

In Chapter 4, we were guided by the question “who is teaching our students?” to catechize the issue of absent pedagogical content knowledge of graduate level interdisciplinary science instructors. To alleviate this issue, we developed and operationalized a pedagogical professional development program specifically for interdisciplinary scientists called the Pedagogy of Interdisciplinary Science Education (POISE) training program.

POISE was made possible through funding and coaching provided by the UMMS RISE program. Additional support came from the insight and guidance of DCMB, our partners, and the POISE advisory council. To continue the study and to broaden the impact of pedagogical training of future interdisciplinary science educators, the first step would be to identify ownership by determining which program, school, or program POISE belongs and answers to. Answering this, the second step would be to identify funding sources. We believe that additional internal iterations of POISE could provide enough background and institutional buy in to make POISE eligible for a research education grant (R25). If so, the funding could support a larger staff and potentially incentivize training through the development of a recognized certification program as well as other professional development resources and training. Last, the initial curriculum of POISE requires additional iterations and insight from the PAC to create a more comprehensive, exhaustive, and inclusive scope and sequence for the trainees.

5.4 Closing

The combined works of this dissertation were unintentionally discovered. I did not expect this journey or this output. I observed a perceived gap in graduate education and my advisors, department, school, and university supported my initiative and curiosity in such a way that this body of work sets a definite precedent of the academic innovation alive at the University of

Michigan. This dissertation could only be described as “alternative academic output,” and while it was not without administrative challenge, it was only possible with the support of the institution. Therefore, this precedent makes it possible for faculty to accept and encourage more diverse approaches to studying bioinformatics as a whole and allows students the flexibility to pursue non-traditional professional development and inquiry without the fear of “staying on track.” This consequence is probably the most impactful future direction of this dissertation, and I am excited to see whether and how it grows.

Appendices

Appendix A Bioinformatics Curriculum Survey

Start of Block: Demographics

It is important to be upfront with the purpose of this survey.

What this survey is and is for:

The purpose of this survey is to ask for insight into identifying potential learning obstacles of bioinformatics—in general—*within our program*. By identifying these potential obstacles, we may be able to emphasize and reorganize the general curriculum of our students to better help them succeed in our program.

Many of the prompts within the survey were provided by student focus groups and faculty advisory boards. This does not mean that the points made in this survey are the *only* points/topics/concepts that are of interest. If, at any point, you would like to provide additional topics/subjects to be added to the research; feel free to use the provided prompts to make those additions.

The information collected through this survey will be anonymized and aggregated. However, be aware that some of the prompts are open-ended. Therefore, respondents should exercise caution when completing these prompts to remove any identifying information should they choose to remain anonymous.

What this survey is not:

This survey is nothing without your input. Please take the time to carefully consider your responses since that data may be used to help students of our program in the future.

The prompts in this survey are not comprehensive and exhaustive. That is, we do not have a *complete* picture of the learning obstacles of our students. Just a starting point.

This survey is not meant as an outlet to give negative/positive feedback about any specific instructor or course

This survey (and collected data) may not change anything in our department in your observable future, but it is important nonetheless.

Disclaimer:

This survey is part of an ongoing research project called the "Identification of threshold concepts within bioinformatics" (HUM00185545).

Is the Department of Computational Medicine and Bioinformatics your primary department?

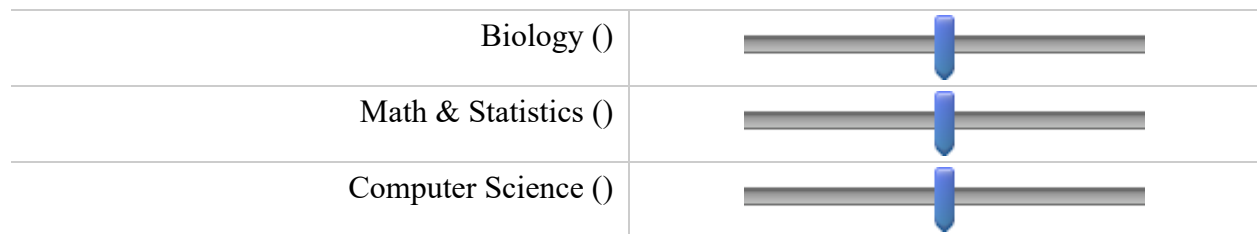
▼ No (1) ... Yes (2)

How long have you been part of the bioinformatics department?

▼ < 1 year (1) ... 4+ years (5)

How would you rate yourself within the following skill sets (1 being the worst)?

1 2 3 4 5 6 6 7 8 9 10



If you had to rank the importance of biology, math & statistics, and computer science to performing your bioinformatic research; what order would they be in (1 being most important)?

_____ Biology (1)

_____ Math & Statistics (2)

_____ Computer Science (3)

What primary sub-field of bioinformatics is either your current or prospective research in?

▼ Computational Biology (1) ... Other (16)

Display This Question:

If What primary sub-field of bioinformatics is either your current or prospective research in? = Other

If "Other," please describe your sub-field:

Do you have any programming experience?

▼ No (1) ... Yes (2)

Display This Question:

If Do you have any programming experience? = Yes

Please select your level of agreement with the statement "I am proficient at programming"

- Strongly Disagree (1)
- Somewhat disagree (2)
- Neither agree nor disagree (3)
- Somewhat agree (4)
- Strongly agree (5)

Display This Question:

If Do you have any programming experience? = Yes

What programming languages do your primarily use?

- Python (1)
- R (2)
- C/C++/C# (3)
- Go (4)
- Rust (5)
- JavaScript (6)
- Shell & Bash (7)
- Excel (8)
- julia (9)
- MATLAB (10)
- Other (11)

Display This Question:

If What programming languages do your primarily use? = Other

If "Other," what other programming language do you use that were not listed?

End of Block: Demographics

Start of Block: Biological Concepts

In this next section, the survey will be focusing on some *biologically specific* concepts relevant to bioinformatics. However, be aware that there are few overlapping and generalized biological concepts that span all sub-disciplines of bioinformatics. This means that a difficult biological concept of one sub-discipline (Genomics) may not even apply to another (Pharmacology). As such, this section may seem overly generalized and/or vague.

At the end of this section, you will have an opportunity to provide us with any information/concepts you feel were missing from this section.

What is your level agreement for the following statements:

	Strongly Disagree (1)	Somewhat disagree (2)	Neither agree nor disagree (3)	Somewhat agree (4)	Strongly agree (5)
I understand the "Central Dogma" of molecular biology (1)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I understand how to perform bioinformatic literature review in my field (2)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I understand the concept and application of references of extant knowledge (e.g. reference genomes, reference sets, and standard curves) (3)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Display This Question:

If What is your level agreement for the following statements: [I understand the "Central Dogma" of molecular biology] (Recode) >= 4

In your own words; please describe the "Central Dogma"

Are there any *biologically-relevant* bioinformatics concepts you have had difficulties understanding that were not listed above?

No (1)

Yes (2)

Display This Question:

If Are there any biologically-relevant bioinformatics concepts you have had difficulties understandi... = Yes

If "Yes," please describe any difficult *biologically-relevant* bioinformatics concept(s) not listed previously

End of Block: Biological Concepts

Start of Block: Statistical Concepts

In this next section, the survey will be focusing on *math and statistics-specific* concepts relevant to bioinformatics. Again, there are very few mathematical concepts that span all (or many) sub-disciplines of bioinformatics. This means that the concepts within this section are somewhat sparse, vague, and/or generalized.

At the end of this section, you will have an opportunity to provide us with any information/concepts you feel were missing from this section.

What level of agreement do you have with the following statements:

	Strongly Disagree (1)	Somewhat disagree (2)	Neither agree nor disagree (3)	Somewhat agree (4)	Strongly agree (5)
When necessary, I know how to interrogate bioinformatic data to determine their statistical properties (6)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
When necessary, I know how to identify the statistical distribution(s) that apply to a given set of bioinformatics data (7)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I understand the differences between "local" and "global" when considering concepts like sequence alignment (8)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Display This Question:

If What level of agreement do you have with the following statements: [I understand the differences between "local" and "global" when considering concepts like sequence alignment] (Recode) >= 4

How would you describe the difference between "local" and "global" as it relates to alignments?

Are there any mathematical or statistically-relevant bioinformatics concepts you have had difficulties understanding that were not covered?

- No (1)
- Yes (2)

Display This Question:

If Are there any mathematical or statistically-relevant bioinformatics concepts you have had difficu... = Yes

If "Yes," please describe any difficult mathematical or statistically-relevant bioinformatics concept(s) not covered

End of Block: Statistical Concepts

Start of Block: Computer Science Concepts

In this next section, the survey will be focusing on computer science-specific concepts relevant to bioinformatics.

At the end of this section, you will have an opportunity to provide us with any information/concepts you feel were missing from this section.

What is your level agreement for the following statements:

	Strongly disagree (1)	Somewhat disagree (2)	Neither agree nor disagree (3)	Somewhat agree (4)	Strongly agree (5)
I understand how to perform a proper code review (1)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I can confidently exercise debugging strategies (2)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I know how to package and deploy a program in a scalable and maintainable way (3)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I can confidently use version control (4)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I can confidently navigate and manipulate the command line interface (5)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I can quickly deconstruct a programmatic task into its constitutive components (6)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Display This Question:

If What is your level agreement for the following statements: [I can confidently use version control] (Recode) >= 4

How do you manage version control?

Display This Question:

If What is your level agreement for the following statements: [I can confidently navigate and manipulate the command line interface] (Recode) >= 3

What is the difference between relative and absolute paths?

What is your level agreement for the following statements:

	Strongly disagree (1)	Somewhat disagree (2)	Neither agree nor disagree (3)	Somewhat agree (4)	Strongly agree (5)
I understand the difference between "lazy" and "eager" loading with respect to data analysis (1)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I understand when and why certain data structures (e.g. dictionaries and dataframes) are used for data analysis (2)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I know how to access and leverage multiple indexing strategies (the representation of an item's position within a sequence) across multiple data structures (3)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I comprehend the concepts, implementation, and application of dynamic programming algorithms (4)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

What is your level agreement for the following statements:

	Strongly disagree (1)	Somewhat disagree (2)	Neither agree nor disagree (3)	Somewhat agree (4)	Strongly agree (5)
I know how to scale prototype algorithms to high-performance and high-throughput computing (1)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I comprehend functional motif searching and identification (2)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I can confidently compute sequence similarities of two (2) or more sequences (4)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Display This Question:

If What is your level agreement for the following statements: [I know how to scale prototype algorithms to high-performance and high-throughput computing] (Recode) >= 4

Please describe one reason why requesting 7 cores for a Python job may not be efficient?

Display This Question:

If What is your level agreement for the following statements: [I comprehend functional motif searching and identification] (Recode) <= 3

Do you know what a "motif" is?

- No (1)
 - Maybe (2)
 - Yes (3)
-

Are there any computer science-specific bioinformatics concepts you have had difficulties understanding that were not covered?

- No (1)
- Yes (2)

Display This Question:

If Are there any computer science-specific bioinformatics concepts you have had difficulties underst... = Yes

If "Yes," please describe any difficult computer science-specific bioinformatics concept(s) not covered

End of Block: Computer Science Concepts

Start of Block: Open-ended

Are there any difficult concepts/topics not covered in the survey that should have been

- Yes (1)
- No (2)

Display This Question:

If Are there any difficult concepts/topics not covered in the survey that should have been = Yes

If "Yes," please describe

Is there anything else you would like us to know that is relevant to identifying obstacles in the bioinformatic curriculum?

- No (1)
- Yes (2)

Display This Question:

If Is there anything else you would like us to know that is relevant to identifying obstacles in the... = Yes

If "Yes," please describe

End of Block: Open-ended

Appendix B DCMB Threshold Concepts Faculty Survey

Start of Block: Preface

Preface **Who is this survey for?**

Faculty associated with the Department of Computational Medicine and Bioinformatics

What this survey is and is for:

This survey is part of my dissertation research into the education of bioinformatics.

The purpose of this survey is to ask for insight into identifying potential learning obstacles of bioinformatics—in general—within our bioinformatics program. By identifying these potential obstacles, we may be able to emphasize and reorganize the general curriculum of our students to better help them succeed in our program

Many of the prompts within the survey were distilled from student focus groups and a survey of the students ($n = 70$). This does not mean that the concepts covered are the only concepts that are of interest. If, at any point, you would like to provide additional concepts to be added to the research; feel free to use the provided prompts to make your suggestions.

The information collected through this survey will be anonymized and aggregated. However, be aware that some of the prompts are open-ended. Therefore, respondents should exercise caution when completing these prompts to remove any identifying information should they choose to remain anonymous.

What this survey is not:

This survey is nothing without your input

Please take the time to carefully consider your responses since that data may be used to help students of our program in the future The prompts in this survey are not comprehensive and exhaustive. That is, we do not have a complete picture of the learning obstacles of our students.

Just a starting point.

This survey is not meant as an outlet to give negative/positive feedback about any specific instructor, student, or course

This survey (and collected data) may not change anything in our department in your observable future, but it is important nonetheless

Disclaimer:

This survey is part of an ongoing research project called the "Identification of threshold concepts within bioinformatics" (HUM00185545).

End of Block: Preface

Start of Block: Demographics

Association What is your association with Department of Computational Medicine and Bioinformatics?

▼ Primary (1) ... CCMB Affiliate (4)

Field of Study What sub-field of bioinformatics best describes your field of research?

▼ Computational Biology (1) ... Other (16)

Teaching Please select your level of agreement to the following statement: I consider 15% or more of my effort to be dedicated to didactic teaching

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Estimate your number of classroom contact hours per year with students

- 0-10 (1)
 - 10-20 (2)
 - 20+ (3)
-

Please estimate how long you have held your current teaching role at the University

Pedagogical Training Do you have any formal pedagogical training? For example; seminars on lesson planning, inclusive teaching, or preparing future faculty workshops.

- Yes, a lot (1)
- Yes, but not much (2)
- No (3)

End of Block: Demographics

Start of Block: Threshold Concept Definition

Definition The main purpose of this survey is to attempt to identify potential Threshold Concepts (TCs) of bioinformatics. TCs are domain-specific concepts that are difficult to understand but nearly impossible to unlearn and fundamentally change the way the learner perceives the domain. Every domain can be defined by its TCs. For example, *Object-Oriented Programming* is a difficult to understand concept that—once understood—cannot be unlearned and ultimately changes the way a learner perceives computer science.

It is important to note that TCs are not “competencies.” Competencies are usually generalized skillsets (e.g. critical thinking or communication) that are gained over time and not through direct education. A TC is a specific concept that can be directly taught during a regular course of study (e.g. gravity).

The approach this study has taken in identifying bioinformatic TCs is by first identifying troublesome concepts through direct student survey and focus groups. These concepts are then broken out by student demographic data (e.g. time in program). Ideally, any concept that is shown to be difficult to understand by junior students but taken for granted by senior students suggests the possibility of a TC.

These potential TCs are what this survey will focus on. You, as the instructors of the bioinformatic field at UM, are asked to disprove/clarify/corroborate these findings. A faculty/staff TC focus group will be brought together at a later date to discuss, refine, and/or eliminate any of these TCs.

End of Block: Threshold Concept Definition

Start of Block: Sequencing

Sequential Data Analysis:

This is a *generalized* term for a recurring concept brought up by students.

The concept of "sequential data analysis" (as defined here) is how the component parts of a dataset are collected/constructed and analyzed. This concept encompasses all sequential analysis of data in which order matters. Examples of sequential data are nucleic acid sequencing, signal processing, machine learning filters/windows, etc.

Please select your level of agreement to the following statement:

Understanding the concept of **sequential data analysis** (as defined above) fundamentally changes how one perceives the field of bioinformatics

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Please select your level of agreement to the following statement:

The concept of **sequential data analysis** (as defined above) is a concept that students can have considerable difficulty understanding

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Please select your level of agreement to the following statement:

The concept of **sequential data analysis** (as defined above) can be or is unique to the field of bioinformatics

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Please select your level of agreement to the following statement:

It is important that students within bioinformatics understand the concept of **sequential data analysis** (as defined above)

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Would you suggest a different description/definition of this concept or would like to refine what has been suggested?

- Yes (1)
- No (2)

Display This Question:

If Would you suggest a different description/definition of this concept or would like to refine what... = Yes

Please use this space to suggest a different description/definition of this concept or refinement:

Do you have any anecdotes or examples from your teaching/mentoring that corroborates this identified concept?

Yes (1)

No (2)

Display This Question:

If Do you have any anecdotes or examples from your teaching/mentoring that corroborates this identif... = Yes

Please use this space to share any anecdotes or examples from your teaching/mentoring that corroborates this identified concept:

Would you like to provide anything else regarding this concept?

Yes (1)

No (2)

Display This Question:

If Would you like to provide anything else regarding this concept? = Yes

Please use this space to provide anything else you would like us to know regarding this concept:

End of Block: Sequencing

Start of Block: Statistics

Statistical distribution(s) identification and application:

This concept was identified by how important it was weighted by senior students and how little the concept was self-reportedly understood.

Please select your level of agreement to the following statement:

Understanding the concept of **statistical distribution(s) identification and application** (as defined above) fundamentally changes how one perceives the field of bioinformatics

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Please select your level of agreement to the following statement:

The concept of **statistical distribution(s) identification and application** (as defined above) is a concept that students can have considerable difficulty understanding

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Please select your level of agreement to the following statement:

The concept of **statistical distribution(s) identification and application** (as defined above) can be or is unique to the field of bioinformatics

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Please select your level of agreement to the following statement:

It is important that students within bioinformatics understand the concept of **statistical distribution(s) identification and application** (as defined above)

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Would you suggest a different description/definition of this concept or would like to refine what has been suggested?

- Yes (1)
- No (2)

Display This Question:

If Would you suggest a different description/definition of this concept or would like to refine what... = Yes

Please use this space to suggest a different description/definition of this concept or refinement:

Do you have any anecdotes or examples from your teaching/mentoring that corroborates this identified concept?

Yes (1)

No (2)

Display This Question:

If Do you have any anecdotes or examples from your teaching/mentoring that corroborates this identif... = Yes

Please use this space to share any anecdotes or examples from your teaching/mentoring that corroborates this identified concept:

Would you like to add anything else regarding this concept?

Yes (1)

No (2)

Display This Question:

If Would you like to add anything else regarding this concept? = Yes

Please use this space to provide anything else you would like us to know regarding this concept:

End of Block: Statistics

Start of Block: Data management

Data ingest, exploration, and management:

This concept focuses on:

1. identifying potential edge cases within a given dataset
2. the principles and approaches to data cleansing (cleaning and normalization)
3. data wrangling: transforming, reformatting, and/or the remapping of data
4. data exploration

Please select your level of agreement to the following statement:

Understanding the concepts of **data ingest, exploration, and management** (as defined above) fundamentally changes how one perceives the field of bioinformatics

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Please select your level of agreement to the following statement:

The concepts of **data ingest, exploration, and management** (as defined above) is a concept that students can have considerable difficulty understanding

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Please select your level of agreement to the following statement:

The concepts of **data ingest, exploration, and management** (as defined above) can be or is unique to the field of bioinformatics

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Please select your level of agreement to the following statement:

It is important that students within bioinformatics understand the concepts of **data ingest, exploration, and management** (as defined above)

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Would you suggest a different description/definition of this concept or would like to refine what has been suggested?

- Yes (1)
- No (2)

Display This Question:

If Would you suggest a different description/definition of this concept or would like to refine what... = Yes

Please use this space to suggest a different description/definition of this concept or refinement:

Do you have any anecdotes or examples from your teaching/mentoring that corroborates this identified concept?

Yes (1)

No (2)

Display This Question:

If Do you have any anecdotes or examples from your teaching/mentoring that corroborates this identif... = Yes

Please use this space to share any anecdotes or examples from your teaching/mentoring that corroborates this identified concept:

Would you like to add anything else regarding this concept?

Yes (1)

No (2)

Display This Question:

If Would you like to add anything else regarding this concept? = Yes

Please use this space to provide anything else you would like us to know regarding this concept:

End of Block: Data management

Start of Block: Standards

References to extant knowledge:

This *generalized* term encompasses ideas such as reference genomes, reference sets, and calibration/standard curves. More specifically how they are made, why they are important, and how are they applied.

Please select your level of agreement to the following statement:
Understanding the concept of **references to extant knowledge** (as defined above)
fundamentally changes how one perceives the field of bioinformatics

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Please select your level of agreement to the following statement:
The concept of **references to extant knowledge** (as defined above) is a concept that students
can have considerable difficulty understanding

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Please select your level of agreement to the following statement:

The concept of **references to extant knowledge** (as defined above) can be or is unique to the field of bioinformatics

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Please select your level of agreement to the following statement:

It is important that students within bioinformatics understand the concept of **references to extant knowledge** (as defined above)

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Would you suggest a different description/definition of this concept or would like to refine what has been suggested?

- Yes (1)
- No (2)

Display This Question:

If Would you suggest a different description/definition of this concept or would like to refine what... = Yes

Please use this space to suggest a different description/definition of this concept or refinement:

Do you have any anecdotes or examples from your teaching/mentoring that corroborates this identified concept?

Yes (1)

No (2)

Display This Question:

If Do you have any anecdotes or examples from your teaching/mentoring that corroborates this identif... = Yes

Please use this space to share any anecdotes or examples from your teaching/mentoring that corroborates this identified concept:

Would you like to add anything else regarding this concept?

Yes (1)

No (2)

Display This Question:

If Would you like to add anything else regarding this concept? = Yes

Please use this space to provide anything else you would like us to know regarding this concept:

End of Block: Standards

Start of Block: Data Scaling

Data Scaling:

Data scaling (as defined here) are the principles, applications, and etiquette of moving an algorithm or analysis pipeline from prototype test cases to high-throughput analysis on high performance computing environments.

Please select your level of agreement to the following statement:

Understanding the concept of **data scaling** (as defined above) fundamentally changes how one perceives the field of bioinformatics

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Please select your level of agreement to the following statement:

The concept of **data scaling** (as defined above) is a concept that students can have considerable difficulty understanding

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Please select your level of agreement to the following statement:

The concept of **data scaling** (as defined above) can be or is unique to the field of bioinformatics

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Please select your level of agreement to the following statement:

It is important that students within bioinformatics understand the concept of **data scaling** (as defined above)

- Strongly disagree (1)
 - Somewhat disagree (2)
 - Neither agree nor disagree (3)
 - Somewhat agree (4)
 - Strongly agree (5)
-

Would you suggest a different description/definition of this concept or would like to refine what has been suggested?

- Yes (1)
- No (2)

Display This Question:

If Would you suggest a different description/definition of this concept or would like to refine what... = Yes

Please use this space to suggest a different description/definition of this concept or refinement:

Do you have any anecdotes or examples from your teaching/mentoring that corroborates this identified concept?

Yes (1)

No (2)

Display This Question:

If Do you have any anecdotes or examples from your teaching/mentoring that corroborates this identif... = Yes

Q53 Please use this space to share any anecdotes or examples from your teaching/mentoring that corroborates this identified concept:

Would you like to add anything else regarding this concept?

Yes (1)

No (2)

Display This Question:

If Would you like to add anything else regarding this concept? = Yes

Please use this space to provide anything else you would like us to know regarding this concept

End of Block: Data Scaling

Start of Block: Block 8

Is there anything else you would like to provide us with or say?

Yes (1)

No (2)

Display This Question:

If Is there anything else you would like to provide us with or say? = Yes

Q63 Please feel free to use the space below to share with us anything else you would like us to know:

End of Block: Block 8

Appendix C POISE Pilot Cohort Feedback

Start of Block: Demographics

What is your year in graduate school?

- 1 (1)
- 2 (2)
- 3 (3)
- 4+ (4)
-

What is your discipline/field of study? (e.g. Biostatistics)

How much prior teaching experience did you have prior to enrolling in POISE?

- Less than 1 term (1)
- 1-2 terms (2)
- 3-4 terms (3)
- 5+ terms (4)

End of Block: Demographics

Start of Block: Post-then-Pre

For each of the statements listed below, please indicate your level of knowledge or perspective both before and after completing the POISE program

	Before POISE	After POISE

	Strongly Disagree (1)	Disagree (2)	Neither Agree or Disagree (3)	Agree (4)	Strongly Agree (5)	Strongly Disagree (1)	Disagree (2)	Neither Agree or Disagree (3)	Agree (4)	Strongly Agree (5)
I have a comprehensive understanding of learning theories (1)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I know how to apply learning theories to develop curriculum (e.g. Bloom's Taxonomy) (2)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I know how to develop effective learning objectives (3)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I know how to use technology to manage a virtual classroom (4)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I know how to use backward design to develop curricula (5)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I understand the principles of inclusive teaching (6)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I know how to create authentic assessments (7)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I feel prepared to teach at the post-secondary level (8)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

End of Block: Post-then-Pre

Start of Block: Multiple-choice

Please indicate your agreement with the following statements:

	Strongly Disagree (1)	Somewhat disagree (2)	Neither agree nor disagree (3)	Somewhat agree (4)	Strongly agree (5)
The POISE learning objectives were clear (1)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
The length of the sessions were appropriate (2)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
The readings & resources were appropriate for the topics (3)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
My participation in POISE changed my teaching confidence (4)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
POISE provided a proportional balance between theoretical knowledge and applied learning (5)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
The topics presented were relevant to my future educator role (6)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
The POISE program met my expectations (7)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
I would recommend POISE to a colleague that is interested in teaching (8)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Please indicate your level of satisfaction with the following:

	Very dissatisfied (1)	Dissatisfied (2)	Neutral (3)	Satisfied (4)	Very Satisfied (5)
Order of modules and content (1)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Each instructor's knowledge about the topics (2)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Use of Zoom for session meetings (3)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
How the sessions were scheduled (4)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Work load expectations for the capstone project (5)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Regarding the topics covered in POISE; please rank them in order of personal impact on your teaching confidence

- _____ Critical self-reflection (1)
- _____ Cognitive development theory & Sociocultural theory (2)
- _____ Bloom's taxonomy (3)
- _____ Cognitive load theory (4)
- _____ Threshold concepts & competency-based education (5)
- _____ Evaluation & assessment (6)
- _____ Virtual classroom management (7)
- _____ Backwards design (8)
- _____ Rubric design (9)
- _____ Inclusive teaching (10)
- _____ Canvas Course administration & technology in the classroom (11)
- _____ Capstone project (12)

End of Block: Multiple-choice

Start of Block: Open-ended Questions

Reflect on your SMART goals that you developed at the beginning of POISE. Which SMART goal(s) did you meet (if any) as a result of your participation in POISE?

How will you apply what you learned through POISE in your future educator role?

What topics or activities were most useful to you in developing as a future educator and what made them useful?

What topics or activities were least useful to you in developing as a future educator and how you would suggest they be changed or improved?

Please share anything else you would like us to know either about POISE, in general or as it pertains to your teaching confidence, that was not asked in this evaluation?

End of Block: Open-ended Questions

Bibliography

- 1 Harrison S, Grant C. Exploring of new models of research pedagogy: time to let go of master-apprentice style supervision? *Teaching in Higher Education* 2015;**20**:556–66. <https://doi.org/10.1080/13562517.2015.1036732>.
- 2 Ryan K. The Transition from Undergraduate to Graduate Student. *Careers in Food Science: From Undergraduate to Professional* 2008:105–16. https://doi.org/10.1007/978-0-387-77391-9_11.
- 3 Xu D, Solanki S. Tenure-Track Appointment for Teaching-Oriented Faculty? The Impact of Teaching and Research Faculty on Student Outcomes: *Educational Evaluation and Policy Analysis* 2019;**42**:66–86. <https://doi.org/10.3102/0162373719882706>.
- 4 Flaherty C. *How many is too many teaching waivers at a public research institution?*. Inside Higher Ed. 2015. URL: <https://www.insidehighered.com/news/2015/12/08/how-many-too-many-teaching-waivers-public-research-institution> (Accessed May 8, 2022).
- 5 Prahalad CK, Hamel G. The core competence of the corporation. *Strategische Unternehmensplanung/Strategische Unternehmensführung*. Springer; 1997. p. 969–87.
- 6 PIBS Curriculum Committee. *Program in the Biomedical Sciences Competency Areas*. Ann Arbor; 2019.
- 7 Slack N. The Importance-Performance Matrix as a Determinant of Improvement Priority. *International Journal of Operations & Production Management* 1994;**14**:59–75. <https://doi.org/10.1108/01443579410056803>.
- 8 Martilla JA, James JC. Importance-Performance Analysis. *Journal of Marketing* 1977;**41**:77–9. <https://doi.org/10.1177/002224297704100112>.
- 9 Dee TS, Jacob B. The impact of no Child Left Behind on student achievement. *Journal of Policy Analysis and Management* 2011;**30**:418–46. <https://doi.org/10.1002/PAM.20586>.
- 10 Darrow A-A. The Every Student Succeeds Act (ESSA). *General Music Today* 2016;**30**:41–4. <https://doi.org/10.1177/1048371316658327>.
- 11 U.S. Department of Education. *Every Student Succeeds Act (ESSA)* . n.d. URL: <https://www.ed.gov/essa?src=rn> (Accessed May 11, 2022).
- 12 Phelps RP. *Teach to the Test?*. The Wilson Quarterly. 2011. URL: <https://www.jstor.org/stable/41484371?seq=1> (Accessed May 11, 2022).
- 13 Ro J. Learning to teach in the era of test-based accountability: a review of research. *Professional Development in Education* 2018;**45**:87–101. <https://doi.org/10.1080/19415257.2018.1514525>.
- 14 Stallings WM, Leslie EK. Student Attitudes toward Grades and Grading. *Improving College and University Teaching*, 1970;**18**:66–8. <https://doi.org/10.1080/00193089.1970.10532929>.
- 15 Dulger M, Deniz H. Assessing the Validity of Multiple-choice Questions in Measuring Fourth Graders' Ability to Interpret Graphs about Motion and Temperature. *INTERNATIONAL JOURNAL OF ENVIRONMENTAL & SCIENCE EDUCATION* 2016;**12**:177–93.

- 16 Jensen JL, McDaniel MA, Woodard SM, Kummer TA. Teaching to the Test...or Testing to Teach: Exams Requiring Higher Order Thinking Skills Encourage Greater Conceptual Understanding. *Educational Psychology Review* 2014;**26**:307–29. <https://doi.org/10.1007/S10648-013-9248-9>.
- 17 McKenna P. Multiple choice questions: answering correctly and knowing the answer. *Interactive Technology and Smart Education* 2019;**16**:59–73. <https://doi.org/10.1108/ITSE-09-2018-0071/FULL/XML>.
- 18 Johnson DT, Price JE. *Will this be on the test? : what your professors really want you to know about succeeding in college*. Princeton University Press; 2019.
- 19 Chamberlin K, Yasué M, Chiang ICA. The impact of grades on student motivation. *Active Learning in Higher Education* 2018. <https://doi.org/10.1177/1469787418819728>.
- 20 Hall JD, O’Connell AB, Cook JG. Predictors of student productivity in biomedical graduate school applications. *PLoS ONE* 2017;**12**:. <https://doi.org/10.1371/JOURNAL.PONE.0169121>.
- 21 Raskin JD. *The Meaning of Grades in Graduate School: Is a “B” in grad school really a “C”?* Psychology Today. 2017. URL: <https://www.psychologytoday.com/us/blog/making-meaning/201701/the-meaning-grades-in-graduate-school> (Accessed May 11, 2022).
- 22 Kuo M. *MBA Course Grading Curve at UCLA Anderson*. MBA Excel. 2012. URL: <https://www.mbaexcel.com/mba/mba-course-grading-curve-at-ucla-anderson/> (Accessed May 11, 2022).
- 23 Cushman T. Who Best to Tame Grade Inflation? *Academic Questions* 2003;**16**:48–56.
- 24 Jamieson JP, Harkins SG. The effect of stereotype threat on the solving of quantitative GRE problems: A mere effort interpretation. *Personality and Social Psychology Bulletin* 2009;**35**:1301–14. <https://doi.org/10.1177/0146167209335165>.
- 25 Astin AW, Antonio AL. *Assessment for Excellence: The Philosophy and Practice of Assessment and Evaluation in Higher Education*. Second Edition. Lanham: Rowman & Littlefield Publishers, Inc.; 2012.
- 26 Walsh B. *When Testing Takes Over* . Harvard Graduate School of Education. 2017. URL: <https://www.gse.harvard.edu/news/uk/17/11/when-testing-takes-over> (Accessed May 11, 2022).
- 27 Petersen SL, Erenrich ES, Levine DL, Vigoreaux J, Gile K. Multi-institutional study of GRE scores as predictors of STEM PhD degree completion: GRE gets a low mark. *PLoS ONE* 2018;**13**:. <https://doi.org/10.1371/JOURNAL.PONE.0206570>.
- 28 Moneta-Koehler L, Brown AM, Petrie KA, Evans BJ, Chalkley R. The Limitations of the GRE in predicting success in biomedical graduate school. *PLoS ONE* 2017;**12**:. <https://doi.org/10.1371/JOURNAL.PONE.0166742>.
- 29 Wiggins GP, Wiggins G, McTighe J. *Understanding by design*. Ascd; 2005.
- 30 Bloom BS. *Taxonomy of educational objectives, handbook I: The cognitive domain*. New York: David McKay Co Inc.; 1956.
- 31 Jackson CK. What do test scores miss? The importance of teacher effects on non–test score outcomes. *Journal of Political Economy* 2018;**126**:2072–107. <https://doi.org/10.1086/699018>.
- 32 Howe C, Hennessy S, Mercer N, Vrikki M, Wheatley L. Teacher-Student Dialogue During Classroom Teaching: Does It Really Impact on Student Outcomes? *Journal of the Learning Sciences* 2019;**28**:462–512. <https://doi.org/10.1080/10508406.2019.1573730>.

- 33 Hanson JM, Paulsen MB, Pascarella ET. Understanding graduate school aspirations: the effect of good teaching practices. *Higher Education* 2016;**71**:735–52. <https://doi.org/10.1007/S10734-015-9934-2>.
- 34 Miller A, Gore J, Wallington C, Harris J, Prieto-Rodriguez E, Smith M. Improving student outcomes through professional development: Protocol for a cluster randomised controlled trial of Quality Teaching Rounds. *International Journal of Educational Research* 2019;**98**:146–58. <https://doi.org/10.1016/J.IJER.2019.09.002>.
- 35 Ebert-May D, Derting TL, Hodder J, Momsen JL, Long TM, Jardeleza SE. What We Say Is Not What We Do: Effective Evaluation of Faculty Professional Development Programs. *BioScience* 2011;**61**:550–8. <https://doi.org/10.1525/BIO.2011.61.7.9>.
- 36 Vilppu H, Södervik I, Postareff L, Murtonen M. The effect of short online pedagogical training on university teachers' interpretations of teaching–learning situations. *Instructional Science* 2019;**47**:679–709. <https://doi.org/10.1007/S11251-019-09496-Z/TABLES/12>.
- 37 Bayram-Jacobs D, Henze I, Evagorou M, Shwartz Y, Aschim EL, Alcaraz-Dominguez S, et al. Science teachers' pedagogical content knowledge development during enactment of socioscientific curriculum materials. *Journal of Research in Science Teaching* 2019;**56**:1207–33. <https://doi.org/https://doi.org/10.1002/tea.21550>.
- 38 Magnusson S, Krajcik J, Borko H. Nature, Sources, and Development of Pedagogical Content Knowledge for Science Teaching. In: Gess-Newsome J, Lederman NG, editors. *Examining Pedagogical Content Knowledge: The Construct and its Implications for Science Education*. Dordrecht: Springer Netherlands; 1999. p. 95–132.
- 39 McCluney CL, Robotham K, Lee S, Smith R, Durkee M. *The Costs of Code-Switching*. Harvard Business Review. 2022. URL: <https://hbr.org/2019/11/the-costs-of-codeswitching> (Accessed May 12, 2022).
- 40 Moore D. Code-switching and Learning in the Classroom. *International Journal of Bilingual Education and Bilingualism* 2002;**5**:279–93. <https://doi.org/10.1080/13670050208667762>.
- 41 Lin A. Classroom code-switching: Three decades of research. *Applied Linguistics Review* 2013;**4**:195–218. <https://doi.org/10.1515/APPLIREV-2013-0009/MACHINEREAADABLECITATION/RIS>.
- 42 Lin AMY. Code-Switching in the Classroom: Research Paradigms and Approaches. *Research Methods in Language and Education* 2017:487–501. https://doi.org/10.1007/978-3-319-02249-9_34.
- 43 Asamoah DA, Doran D, Schiller S. Interdisciplinarity in Data Science Pedagogy: A Foundational Design. *Journal of Computer Information Systems* 2020;**60**:370–7. <https://doi.org/10.1080/08874417.2018.1496803>.
- 44 Wiegard D. *Bioinformatics*. Electrical and Computer Engineering Design Handbook. 2022.
- 45 Chilana PK, Palmer CL, Ko AJ. Comparing Bioinformatics Software Development by Computer Scientists and Biologists: An Exploratory Study. *SECSE* 2009:72–9. <https://doi.org/10.1109/SECSE.2009.5069165>.
- 46 Pathak RK, Singh DB, Singh R. Introduction to basics of bioinformatics. *Bioinformatics* 2022:1–15. <https://doi.org/10.1016/B978-0-323-89775-4.00006-7>.

- 47 *Is Programming Knowledge Necessary For Career In Bioinformatics?*. Biotecnika. 2019. URL: <https://www.biotecnika.org/2019/05/programming-knowledge-necessary-for-career-in-bioinformatics/> (Accessed June 29, 2022).
- 48 Corne DW, Fogel GB. An Introduction to Bioinformatics for Computer Scientists. *Evolutionary Computation in Bioinformatics* 2003;3–18. <https://doi.org/10.1016/B978-155860797-2/50003-2>.
- 49 Dale R, Grüning B, Sjödin A, Rowe J, Chapman BA, Tomkins-Tinch CH, *et al.* Bioconda: Sustainable and comprehensive software distribution for the life sciences. *Nature Methods* 2018;**15**:475–6. <https://doi.org/10.1038/S41592-018-0046-7>.
- 50 Betti E, Cesati M, Gioiosa R, Piermaria F. A global operating system for HPC clusters. *Proceedings - IEEE International Conference on Cluster Computing, ICC* 2009. <https://doi.org/10.1109/CLUSTER.2009.5289191>.
- 51 Chin M. *File Not Found*. The Verge. 2021. URL: <https://www.theverge.com/22684730/students-file-folder-directory-structure-education-gen-z> (Accessed May 12, 2022).
- 52 Nelson C. *Defining Academic Freedom*. Inside Higher Ed. 2010. URL: <https://www.insidehighered.com/views/2010/12/21/defining-academic-freedom> (Accessed May 12, 2022).
- 53 Glover KM. Graduate Student Life Brief. *National Association of Graduate-Professional Students* 2019. <https://doi.org/10.1006/jvbe.2001.1804>.
- 54 Nousiainen MT, Caverzagie KJ, Ferguson PC, Frank JR. Implementing competency-based medical education: What changes in curricular structure and processes are needed? *Medical Teacher* 2017;**39**:594–8. <https://doi.org/10.1080/0142159X.2017.1315077>.
- 55 Loertscher J, Green D, Lewis JE, Lin S, Minderhout V. Identification of Threshold Concepts for Biochemistry. *Https://DoiOrg/101187/Cbe14-04-0066* 2017;**13**:516–28. <https://doi.org/10.1187/CBE.14-04-0066>.
- 56 Nicola-Richmond K, Pépin G, Larkin H, Taylor C. Threshold concepts in higher education: a synthesis of the literature relating to measurement of threshold crossing. *Https://DoiOrg/101080/0729436020171339181* 2017;**37**:101–14. <https://doi.org/10.1080/07294360.2017.1339181>.
- 57 Huq A, Aryal B, Nichols M. “Building blocks: Threshold concepts as interdisciplinary structures of learning” n.d.
- 58 Neve H, Wearn A, Collett T. What are threshold concepts and how can they inform medical education? *Medical Teacher* 2016;**38**:850–3. <https://doi.org/10.3109/0142159X.2015.1112889>.
- 59 Holley KA. The Role of Threshold Concepts in an Interdisciplinary Curriculum: a Case Study in Neuroscience. *Innovative Higher Education* 2017 *43:1* 2017;**43**:17–30. <https://doi.org/10.1007/S10755-017-9408-9>.
- 60 Horrigan LA. Tackling the threshold concepts in physiology: what is the role of the laboratory class? *Adv Physiol Educ* 2018;**42**:507–15. <https://doi.org/10.1152/advan.00123.2017.-Laboratory>.
- 61 Meyer J, Land R. *Threshold concepts and troublesome knowledge: Linkages to ways of thinking and practising within the disciplines*. Citeseer; 2003.
- 62 Delaney SK, Mills J, Galea A, LeBard R, Wilson J, Gibson KJ, *et al.* Analysis of Alternative Strategies for the Teaching of Difficult Threshold Concepts in Large

- Undergraduate Medicine and Science Classes. *Medical Science Educator* 2017;**27**:673–84. <https://doi.org/10.1007/S40670-017-0453-X/TABLES/3>.
- 63 Hyde S, Flatau A, Wilson D. Integrating threshold concepts with reflective practice: Discussing a theory-based approach for curriculum refinement in dental education. *European Journal of Dental Education* 2018;**22**:e687–97. <https://doi.org/10.1111/EJE.12380>.
- 64 Laukkonen RE, Kaveladze BT, Tangen JM, Schooler JW. The dark side of Eureka: Artificially induced Aha moments make facts feel true. *Cognition* 2020;**196**:104122. <https://doi.org/10.1016/J.COGNITION.2019.104122>.
- 65 Sweller J. Cognitive Load Theory. *Psychology of Learning and Motivation - Advances in Research and Theory* 2011;**55**:37–76. <https://doi.org/10.1016/B978-0-12-387691-1.00002-8>.
- 66 Ren X, Kuan PF. Negative binomial additive model for RNA-Seq data analysis. *BMC Bioinformatics* 2020;**21**:1–15. <https://doi.org/10.1186/S12859-020-3506-X/TABLES/1>.
- 67 NARST. *Pedagogical Content Knowledge: Teachers' Integration of Subject Matter, Pedagogy, Students, and Learning Environments*. Research Matters #9702. 1997. URL: <https://narst.org/research-matters/pedagogical-content-knowledge> (Accessed September 29, 2021).
- 68 Henderson C, Beach A, Finkelstein N. Facilitating Change in Undergraduate STEM Instructional Practices: An Analytic Review of the Literature. *J Res Sci Teach* 2011;**48**:952–84. <https://doi.org/10.1002/tea.20439>.
- 69 Nederbragt L. *Scaling Participatory live coding in an undergraduate computational biology course*. In between Lines of Code. 2020. URL: <https://lexnederbragt.com/blog/2020-04-02-scaling-live-coding/> (Accessed May 12, 2022).
- 70 Bhatia A. *Active Learning Leads to Higher Grades and Fewer Failing Students in Science, Math, and Engineering*. WIRED. 2014. URL: <https://www.wired.com/2014/05/empzeal-active-learning/> (Accessed May 12, 2022).
- 71 Steel JJ. Genome Analysis of SARS-CoV-2 Case Study: An Undergraduate Online Learning Activity to Introduce Bioinformatics, BLAST, and the Power of Genome Databases †. *Journal of Microbiology & Biology Education* 2021;**22**:. https://doi.org/10.1128/JMBE.V22I1.2245/SUPPL_FILE/JMBE00023-21_SUPP_1_SEQ2.PDF.
- 72 Waldrop MM. Why we are teaching science wrong, and how to make it right. *Nature* 2015;**523**:272–4. <https://doi.org/10.1038/523272A>.
- 73 Nederbragt LN. *Active learning strategies for bioinformatics teaching | In between lines of code*. In between Lines of Code. 2015. URL: <https://flxlexblog.wordpress.com/2015/08/31/active-learning-strategies-for-bioinformatics-teaching-2/> (Accessed May 12, 2022).
- 74 Freeman S, Eddy SL, McDonough M, Smith MK, Okoroafor N, Jordt H, *et al*. Active learning increases student performance in science, engineering, and mathematics. *Proc Natl Acad Sci U S A* 2014;**111**:8410–5. https://doi.org/10.1073/PNAS.1319030111/SUPPL_FILE/PNAS.1319030111.ST04.DOCX.
- 75 Lee CJ, Toven-Lindsey B, Shapiro C, Soh M, Mazrouee S, Levis-Fitzgerald M, *et al*. Error-discovery learning boosts student engagement and performance, while reducing

- student attrition in a bioinformatics course. *CBE Life Sciences Education* 2018;**17**:.
<https://doi.org/10.1187/CBE.17-04-0061/ASSET/IMAGES/LARGE/CBE-17-AR40-G008.JPEG>.
- 76 Grunspan DZ, Kline MA, Brownell SE. The lecture machine: A cultural evolutionary model of pedagogy in higher education. *CBE Life Sciences Education* 2018;**17**:.
<https://doi.org/10.1187/CBE.17-12-0287>.
- 77 Flaherty C. *Classroom 911: Professor calls the police on two tardy Black students*. Insider Higher Ed. 2022. URL: <https://www.insidehighered.com/news/2022/04/04/professor-calls-police-two-tardy-black-students> (Accessed May 12, 2022).
- 78 D'Angelo F, Iliev N. *Teaching Mathematics to Young Children through the Use of Concrete and Virtual Manipulatives*. 2012. URL: <https://eric.ed.gov/?id=ED534228> (Accessed May 24, 2022).
- 79 Kamina P, Iyer NN. From Concrete to Abstract: Teaching for Transfer of Learning when Using Manipulatives. *NERA Conference Proceedings* 2009;**6**:10–23.
- 80 Tomlinson CA. *What Is Differentiated Instruction?*. Reading Rockets. 2009. URL: <https://www.readingrockets.org/article/what-differentiated-instruction> (Accessed May 24, 2022).
- 81 Tomlinson CA. The Goals of Differentiation. *Educational Leadership* 2008;**66**:.
82 Tomlinson CA. *The Differentiated Classroom: Responding to the Needs of All Learners*. 2nd ed. Alexandria, VA: ASCD; 2013.
- 83 *Understanding Manipulatives*. SchoolMart. 2017. URL: <https://www.schoolmart.com/2017/10/25/understanding-manipulatives/> (Accessed May 24, 2022).
- 84 Vygotsky LS. *Mind in society: The development of higher psychological processes*. Cambridge, Mass: Harvard University Press; 1987.
- 85 Piaget J. Part I: Cognitive Development in Children--Piaget Development and Learning. *J Res Sci Teach* 2003;**40**:.
86 *A Brief Guide to Genomics*. National Human Genome Research Institute. n.d. URL: <https://www.genome.gov/about-genomics/fact-sheets/A-Brief-Guide-to-Genomics> (Accessed May 29, 2022).
- 87 Feingold EA, Good PJ, Guyer MS, Kamholz S, Liefer L, Wetterstrand K, *et al*. The ENCODE (ENCyclopedia Of DNA Elements) Project. *Science* 2004;**306**:636–40.
<https://doi.org/10.1126/SCIENCE.1105136>.
- 88 Milward EA, Shahandeh A, Heidari M, Johnstone DM, Daneshi N, Hondermarck H. Transcriptomics. *Encyclopedia of Cell Biology* 2016;**4**:160–5.
<https://doi.org/10.1016/B978-0-12-394447-4.40029-5>.
- 89 Besser J, Carleton HA, Gerner-Smidt P, Lindsey RL, Trees E. Next-Generation Sequencing Technologies and their Application to the Study and Control of Bacterial Infections. *Clin Microbiol Infect* 2018;**24**:335. <https://doi.org/10.1016/J.CMI.2017.10.013>.
- 90 Lipman DJ, Pearson WR. Rapid and sensitive protein similarity searches. *Science* 1985;**227**:1435–41. <https://doi.org/10.1126/SCIENCE.2983426>.
- 91 Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biology* 2009;**10**:.
<https://doi.org/10.1186/GB-2009-10-3-R25>.

- 92 Li H, Ruan J, Durbin R. Mapping short DNA sequencing reads and calling variants using mapping quality scores. *Genome Research* 2008;**18**:1851–8. <https://doi.org/10.1101/GR.078212.108>.
- 93 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 2009;**25**:2078–9. <https://doi.org/10.1093/BIOINFORMATICS/BTP352>.
- 94 Cock PJA, Fields CJ, Goto N, Heuer ML, Rice PM. The Sanger FASTQ file format for sequences with quality scores, and the Solexa/Illumina FASTQ variants. *Nucleic Acids Research* 2010;**38**:1767. <https://doi.org/10.1093/NAR/GKP1137>.
- 95 Deutsch LP. *RFC 1952 - GZIP file format specification version 4.3*. Network Working Group. 1996. URL: <https://datatracker.ietf.org/doc/html/rfc1952> (Accessed May 30, 2022).
- 96 Uppsala Multidisciplinary Center for Advanced Computational Science. *Using CRAM to compress BAM files*. Uppsala University, Sweden. n.d. URL: <https://www.uppmax.uu.se/support/user-guides/using-cram-to-compress-bam-files/> (Accessed May 30, 2022).
- 97 Bonfield JK, Marshall J, Danecek P, Li H, Ohan V, Whitwham A, *et al.* HTSlib: C library for reading/writing high-throughput sequencing data. *Gigascience* 2021;**10**:1–6. <https://doi.org/10.1093/GIGASCIENCE/GIAB007>.
- 98 *Pysam is a Python module for reading and manipulating SAM/BAM/VCF/BCF files. It's a lightweight wrapper of the htslib C-API, the same one that powers samtools, bcftools, and tabix*. Github. 2022. URL: <https://github.com/pysam-developers/pysam> (Accessed May 30, 2022).
- 99 Sherman MD, Mills RE. BAMnostic: an OS-agnostic toolkit for genomic sequence analysis. *Journal of Open Source Software* 2018;**3**:826. <https://doi.org/10.21105/JOSS.00826>.
- 100 D'Haeseleer P. What are DNA sequence motifs? *Nature Biotechnology* 2006 **24**:4 2006;**24**:423–5. <https://doi.org/10.1038/nbt0406-423>.
- 101 Schiller MR. Minimotoif Miner: A Computational Tool to Investigate Protein Function, Disease, and Genetic Diversity. *Current Protocols in Protein Science* 2007;**48**:2.12.1-2.12.14. <https://doi.org/10.1002/0471140864.PS0212S48>.
- 102 Weirauch MT, Cote A, Norel R, Annala M, Zhao Y, Riley TR, *et al.* Evaluation of methods for modeling transcription-factor sequence specificity. *Nat Biotechnol* 2013;**31**:126. <https://doi.org/10.1038/NBT.2486>.
- 103 Stormo GD. DNA binding sites: representation and discovery. *Bioinformatics* 2000;**16**:16–23. <https://doi.org/10.1093/BIOINFORMATICS/16.1.16>.
- 104 Altarawy D, Ismail MA, Ghanem SM. MProfiler: A Profile-Based Method for DNA Motif Discovery. *Lecture Notes in Computer Science (Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 2009;**5780 LNBI**:13–23. https://doi.org/10.1007/978-3-642-04031-3_2.
- 105 Crooks GE, Hon G, Chandonia JM, Brenner SE. WebLogo: a sequence logo generator. *Genome Res* 2004;**14**:1188–90. <https://doi.org/10.1101/GR.849004>.
- 106 Bembom O, Ivánek R. *Sequence logos for DNA sequence alignments*. Bioconductor. 2022. URL: <https://bioconductor.org/packages/devel/bioc/vignettes/seqLogo/inst/doc/seqLogo.html> (Accessed May 30, 2022).

- 107 Gentleman RC, Carey VJ, Bates DM, Bolstad B, Dettling M, Dudoit S, *et al.* Bioconductor: open software development for computational biology and bioinformatics. *Genome Biology* 2004 5:10 2004;5:1–16. <https://doi.org/10.1186/GB-2004-5-10-R80>.
- 108 Jang MH. *Linux annoyances for geeks*. O'Reilly; 2006.
- 109 Harris CR, Millman KJ, van der Walt SJ, Gommers R, Virtanen P, Cournapeau D, *et al.* Array programming with NumPy. *Nature* 2020 585:7825 2020;585:357–62. <https://doi.org/10.1038/s41586-020-2649-2>.
- 110 The pandas development team. pandas-dev/pandas: Pandas 2020. <https://doi.org/10.5281/zenodo.3509134>.
- 111 McKinney W. Data Structures for Statistical Computing in Python. *Proceedings of the 9th Python in Science Conference* 2010:56–61. <https://doi.org/10.25080/MAJORA-92BF1922-00A>.
- 112 Tareen A, Kinney JB. Logomaker: beautiful sequence logos in Python. *Bioinformatics* 2020;36:2272–4. <https://doi.org/10.1093/BIOINFORMATICS/BTZ921>.
- 113 Lal A, Ferrarini MG, Gruber AJ. Investigating the human host - ssRNA virus interaction landscape using the SMEAGOL toolbox. *BioRxiv* 2021:2021.12.02.470930. <https://doi.org/10.1101/2021.12.02.470930>.
- 114 van Helden J. Regulatory Sequence Analysis Tools. *Nucleic Acids Research* 2003;31:3593–6. <https://doi.org/10.1093/NAR/GKG567>.
- 115 *Academic Intervention*. APA Dictionary of Psychology. n.d. URL: <https://dictionary.apa.org/academic-intervention> (Accessed June 6, 2022).
- 116 Lynch M. *Types of Classroom Interventions*. The Edvocate. 2019. URL: <https://www.theedadvocate.org/types-of-classroom-interventions/> (Accessed June 6, 2022).
- 117 Charman T, Hepburn S, Lewis M, Lewis M, Steiner A, Rogers SJ, *et al.* Educational Interventions. *Encyclopedia of Autism Spectrum Disorders* 2013:1061–2. https://doi.org/10.1007/978-1-4419-1698-3_1457.
- 118 Dooley Kevin. *Designing large-scale LANs*. O'Reilly; 2002.
- 119 DiFlorio I, Martin B, Middlemiss MA, Duncan PA. Curriculum evaluation. *Nurse Educ Today* 1989;9:402–7. [https://doi.org/10.1016/0260-6917\(89\)90095-6](https://doi.org/10.1016/0260-6917(89)90095-6).
- 120 Walck-Shannon E, Batzli J, Pultorak J, Boehmer H. Biological Variation as a Threshold Concept: Can We Measure Threshold Crossing? *CBE Life Sciences Education* 2019;18:. <https://doi.org/10.1187/CBE.18-12-0241>.
- 121 Tibell LAE, Harms U. Biological Principles and Threshold Concepts for Understanding Natural Selection: Implications for Developing Visualizations as a Pedagogic Tool. *Science and Education* 2017;26:953–73. <https://doi.org/10.1007/S11191-017-9935-X/TABLES/2>.
- 122 Cousin G. An introduction to threshold concepts. *Planet* 2006;17:4–5. <https://doi.org/10.11120/PLAN.2006.00170004>.
- 123 Boustedt J, Eckerdal A, McCartney R, Moström JE, Ratcliffe M, Sanders K, *et al.* Threshold concepts in computer science. *ACM SIGCSE Bulletin* 2007;39:504–8. <https://doi.org/10.1145/1227504.1227482>.
- 124 Dunne T, Low T, Ardington C. Exploring Threshold concepts in basic Statistics, using the Internet. *IASE/ISI Satellite* 2007.
- 125 Meyer JHF, Land R, Baillie C, editors. *Threshold Concepts and Transformational Learning*. vol. 42. Sense Publishers; 2010.

- 126 Moran J. *Interdisciplinarity*. 2nd ed. London and New York: Routledge; 2010.
- 127 *The Integrated Postsecondary Education Data System*. NCES. 2020. URL: <https://nces.ed.gov/ipeds/use-the-data> (Accessed June 8, 2022).
- 128 Glaser BG, Strauss AL. Discovery of grounded theory: Strategies for qualitative research. *Discovery of Grounded Theory: Strategies for Qualitative Research* 2017:1–271. <https://doi.org/10.4324/9780203793206>.
- 129 Tie YC, Birks M, Francis K. Grounded theory research: A design framework for novice researchers. *SAGE Open Medicine* 2019;7:205031211882292. <https://doi.org/10.1177/2050312118822927>.
- 130 Collins H, Callaghan D. What a Difference a Zoom Makes: Intercultural Interactions Between Host and International Students. *Journal of Comparative & International Higher Education* 2022;14:96–111. <https://doi.org/10.32674/JCIHE.V14I2.4300>.
- 131 Qualtrics. Qualtrics 2005.
- 132 Cook DA, Hatala R. Validation of educational assessments: a primer for simulation and beyond. *Advances in Simulation* 2016 1:1 2016;1:1–12. <https://doi.org/10.1186/S41077-016-0033-Y>.
- 133 de Silva MJ, Breuer E, Lee L, Asher L, Chowdhary N, Lund C, *et al*. Theory of Change: A theory-driven approach to enhance the Medical Research Council’s framework for complex interventions. *Trials* 2014;15:1–13. <https://doi.org/10.1186/1745-6215-15-267/FIGURES/2>.
- 134 de Silva MJ, Breuer E, Lee L, Asher L, Chowdhary N, Lund C, *et al*. Theory of Change: A theory-driven approach to enhance the Medical Research Council’s framework for complex interventions. *Trials* 2014;15:1–13. <https://doi.org/10.1186/1745-6215-15-267/FIGURES/2>.
- 135 Kirkpatrick D. Techniques for Evaluating Training Programmes”. *J Am Soc for Train Dev* 1959;11:1–13.
- 136 Cahapay MB. Kirkpatrick Model: Its Limitations as Used in Higher Education Evaluation. *International Journal of Assessment Tools in Education* 2021;8:135–44. <https://doi.org/10.21449/IJATE.856143>.
- 137 *Actionable Theory Of Change (TOC) Guide*. Sopact. 2022. URL: <https://www.sopact.com/theory-of-change> (Accessed June 16, 2022).
- 138 Vygotsky LS. *Mind in Society The Development of Higher Psychological Processes*. Cambridge, Massachusetts: Harvard University Press; 1978.
- 139 Bloom DA, Reid JR, Cassady CI. Education in the time of COVID-19 n.d. <https://doi.org/10.1007/s00247-020-04728-8/Published>.
- 140 Fitzgerald JT, Burkhardt JC, Kasten SJ, Mullan PB, Santen SA, Sheets KJ, *et al*. Assessment challenges in competency-based education: A case study in health professions education. *Medical Teacher* 2015;38:482–90. <https://doi.org/10.3109/0142159X.2015.1047754>.
- 141 Gruppen LD, Burkhardt JC, Fitzgerald JT, Funnell M, Haftel HM, Lypson ML, *et al*. Competency-based education: programme design and challenges to implementation. *Medical Education* 2016;50:532–9. <https://doi.org/10.1111/MEDU.12977>.
- 142 Bibler Zaidi NL, Monrad SU, Grob KL, Gruppen LD, Cherry-Bukowiec JR, Santen SA. Building an Exam Through Rigorous Exam Quality Improvement. *Medical Science Educator* 2017;27:793–8. <https://doi.org/10.1007/S40670-017-0469-2/FIGURES/1>.

- 143 Santen SA, Grob KL, Monrad SU, Stalburg CM, Smith G, Hemphill RR, *et al.* Employing a Root Cause Analysis Process to Improve Examination Quality. *Academic Medicine* 2019;**94**:71–5. <https://doi.org/10.1097/ACM.0000000000002439>.
- 144 Bibler Zaidi NL, Grob KL, Yang J, Santen SA, Monrad SU, Miller JM, *et al.* Theory, Process, and Validation Evidence for a Staff-Driven Medical Education Exam Quality Improvement Process. *Medical Science Educator* 2016;**26**:331–6. <https://doi.org/10.1007/S40670-016-0275-2/TABLES/2>.
- 145 Zaidi NLB, Kreiter CD, Castaneda PR, Schiller JH, Yang J, Grum CM, *et al.* Generalizability of competency assessment scores across and within clerkships: How students, assessors, and clerkships matter. *Academic Medicine* 2018;**93**:1212–7. <https://doi.org/10.1097/ACM.0000000000002262>.
- 146 Zaidi NLB, Grob KL, Monrad SM, Kurtz JB, Tai A, Ahmed AZ, *et al.* Pushing Critical Thinking Skills with Multiple-Choice Questions: Does Bloom’s Taxonomy Work? *Academic Medicine* 2018;**93**:856–9. <https://doi.org/10.1097/ACM.0000000000002087>.
- 147 Davis G, Davis GA. Using A Retrospective Pre-Post Questionnaire To Determine Program Impact. *Annual Meeting of the Mid-Western Educational Research Association* 2002.
- 148 Davis GA. Using a retrospective pre-post questionnaire to determine program impact 2003;**41**..
- 149 Lang D, Savageau JA. Starting at the End: Measuring Learning Using Retrospective Pre-Post Evaluations. *Center for Health Policy and Research (CHPR) Publications* 2017.
- 150 Howard GS. Response-Shift Bias: A Problem in Evaluating Interventions with Pre/Post Self-Reports. *Evaluation Review* 2016;**4**:93–106. <https://doi.org/10.1177/0193841X8000400105>.
- 151 Kirkpatrick JD. *Kirkpatrick’s Four Levels of Training Evaluation*. 1st ed. Association for Talent Development; 2016.
- 152 Heydari MR, Taghva F, Amini M, Delavari S. Using Kirkpatrick’s model to measure the effect of a new teaching and learning methods workshop for health care staff. *BMC Research Notes* 2019;**12**:388. <https://doi.org/10.1186/s13104-019-4421-y>.
- 153 Bates R. A critical analysis of evaluation practice: the Kirkpatrick model and the principle of beneficence. *Evaluation and Program Planning* 2004;**27**:341–7. <https://doi.org/https://doi.org/10.1016/j.evalprogplan.2004.04.011>.
- 154 Wright MC, Finelli CJ, Meizlish D, Bergom I. Facilitating the Scholarship of Teaching and Learning at a Research University. *Change: The Magazine of Higher Learning* 2011;**43**:50–6. <https://doi.org/10.1080/00091383.2011.550255>.
- 155 Schram LN, Pinder-Grover T, Turcic S. Assessing the Long-Term Impact of the Preparing Future Faculty Seminar. *To Improve the Academy* 2017;**36**:101–16. <https://doi.org/10.1002/TIA2.20063>.
- 156 Uno J, Walton KLW. Young Investigator Perspectives. Teaching and the postdoctoral experience: impact on transition to faculty positions. *American Journal of Physiology* 2014;**306**:. <https://doi.org/10.1152/AJPGI.00007.2014>.
- 157 Carroll MA, Catapane EJ, Soto M, Brewer G. NIH/IRACDA Program – a Win for Both Post-Docs and PUI Partner Institutions. *The FASEB Journal* n.d.;**30**:116.3-116.3. https://doi.org/10.1096/FASEBJ.30.1_SUPPLEMENT.116.3.

- 158 Uno JK, Rybarczyk B, Lund PK, Lerea LS, Dykstra L. SPIRE: An Innovative Approach to Postdoctoral Training. *The FASEB Journal* 2009;**23**:632.17-632.17. https://doi.org/10.1096/FASEBJ.23.1_SUPPLEMENT.632.17.
- 159 Broujin de N. A combinatorial problem. *Proceedings of the Section of Sciences of the Koninklijke Nederlandse Akademie van Wetenschappen Te Amsterdam* 1946;**49**:758–64.
- 160 Baum LE, Petrie T. Statistical Inference for Probabilistic Functions of Finite State Markov Chains. *Ann Math Statist* 1966;**37**:1554–63. <https://doi.org/10.1214/AOMS/1177699147>.