# Strategies for Building Computing Skills To Support Microbiome Analysis: a Five-Year Perspective from the EDAMAME Workshop

Ashley Shade,[a,b] Taylor K. Dunivin,[a] Jinlyung Choi,[c] Tracy K. Teal,[d] Adina C. Howe[c]

[a]Department of Microbiology and Molecular Genetics, Michigan State University, East Lansing, Michigan, USA
[b]Department of Plant, Soil and Microbial Sciences, Michigan State University, East Lansing, Michigan, USA
[c]Department of Agricultural and Biosystems Engineering, Iowa State University, Ames, Iowa, USA
[d]Data Carpentry, Davis, California, USA

**ABSTRACT** Here, we report our educational approach and learner evaluations of the first 5 years of the Explorations in Data Analysis for Metagenomic Advances in Microbial Ecology (EDAMAME) workshop, held annually at Michigan State University's Kellogg Biological Station from 2014 to 2018. We hope this information will be useful for others who want to organize computing-intensive workshops and will encourage quantitative skill development among microbiologists.

**IMPORTANCE** High-throughput sequencing and related statistical and bioinformatic analyses have become routine in microbiology in the past decade, but there are few formal training opportunities to develop these skills. A weeklong workshop can offer sufficient time for novices to become introduced to best computing practices and common workflows in sequence analysis. We report our experiences in executing such a workshop targeted to professional learners (graduate students, postdoctoral scientists, faculty, and research staff).

It is now recognized that microbial communities ("microbiomes") play essential roles in the health of the environments and the hosts that they inhabit. In addition, advances in high-throughput sequencing technologies allow observations of the diversity and functional potential of microbiomes in their habitats (1), captured with spatially and temporally ambitious study designs (2). Together, these advances in knowledge and methodology deepen and broaden our understanding of the centrality of microbiomes for host and environmental health. Because of the economy and accessibility of high-throughput sequencing, researchers can now investigate the diversity of interesting microbiomes and can begin to untangle how this diversity contributes to host or ecosystem health. Efforts to capitalize on the promise of microbiome sequencing data have resulted in information-rich genomic data sets that must be analyzed to gain knowledge of their intricate relationships.

We realized that there was a need for broad computational training in microbiome analysis. In 2014, we were encouraged by C. Titus Brown (now at the University of California, Davis) to offer a microbiome analysis workshop. At the time, he led the Analyzing Next-Gen Sequencing (ANGUS; https://angus.readthedocs.io/en/2018/index .html) Workshop at Michigan State University's Kellogg Biological Station (KBS). He noted that some ANGUS learners were particularly interested in microbiome analysis and that there were limited offerings for this training. At the time, there were several

Develop working proficiency at the command line and with shell.

Explain the process of high-throughput sequencing, provide an overview of data-handling (quality control, pre-treatment), and discuss their biases.

Access computing resources: Transfer data and run analyses on Amazon EC2 and/or a high-performance computing cluster.

Access and/or create version-controlled code and resources on GitHub.

Discuss steps in the ecological analyses of microbiomes, including alpha and beta-diversity, ordinations, and resemblance metrics.

Explore datasets and statistically test hypotheses in R.

Visualize patterns in microbial communities using R.

Develop a working proficiency with amplicon sequencing workflows and tools (e.g., QIIME, or mothur, or usearch).

Develop a working proficiency with shotgun metagenomics workflows.

Become familiar with publicly accessible microbial sequence databases/repositories (e.g., NCBI, MG-RAST, FunGene) and the tools that they offer for deposition and analyses.

Identify resources for troubleshooting. This includes: how to ask for and where to find general help online, through peer networks, and from workflow-specific resources (e.g., public tutorials and wikis).

**FIG 1** Overview of learning objectives for the EDAMAME workshop.

short-duration workshops focused on specific tools, such as QIIME (3) and mothur (4), as well as a broader, multiweek course, like STAMPS (https://www.mbl.edu/education/courses/stamps/), at the Marine Biological Laboratory in Woods Hole, MA, USA. There were few workshops that addressed the needs of learners who wanted more information than could be covered in a day but also could not commit to spending several weeks away. Thus, we suspected that there was a need for broad and economical training in microbiome analysis, especially in the U.S. Midwest.

In response, we created a 1-week intensive course to train biologists (from graduate students to faculty) in microbiome-associated sequencing analysis, from raw sequence handling and quality control to statistical analyses and experimental design. We named the course EDAMAME: Explorations in Data Analysis for Metagenomic Advances in Microbial Ecology. Ashley Shade, at the time a new assistant professor in microbial ecology at the Department of Microbiology and Molecular Genetics at Michigan State University, initiated the workshop and started its content development from her materials from a short workshop that she had offered while training in her postdoctoral advisor's lab. Tracy Teal was recruited and brought her array of experience and perspective as a leader in the Software and Data Carpentries workshops, which provide general computing training. In the first year, J. Herr, a postdoc in Shannon Manning's lab at Michigan State who had Data Carpentry training, contributed to developing and implementing the original content. The instruction team expanded in 2016 to include Adina Howe, who was a new faculty member at Iowa State and brought important expertise in untargeted metagenome analysis.

Here, we report a 5-year perspective on the EDAMAME workshop. We describe EDAMAME's learning objectives, target audience and admissions, instructional team, learning environment, educational strategy and assessment, and community resources. We discuss results from assessment, lessons learned, and an outlook for future microbiome training.

## RESULTS

**EDAMAME learning objectives.** EDAMAME's learning objectives were tailored annually to incorporate learners' changing interests and changes in tools and technology (Fig. 1). As a consequence, though the overall content was not drastically changed, we created and retired tutorials as demands changed. For example, when we found that many of our learners had had exposure to and experience with amplicon sequence analysis, we reduced that content to provide more time for metagenome analysis. However, each year featured foundational tutorials in computing literacy, state-of-the-
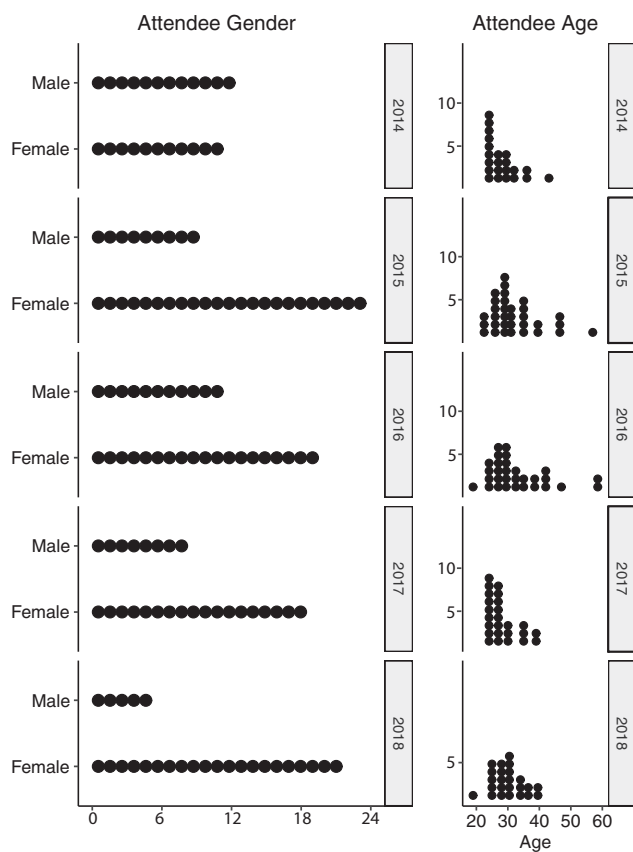
**FIG 2** Distributions of EDAMAME learner sex and age (in years), 2014 to 2018. Data were collected and summarized from pre- and postworkshop assessments. One participant is shown by one dot, and the numbers of participants are on the x axis in the left panel and the y axis in the right panel.

science tools for microbiome analyses, ecological statistics, and computing best practices. We provided specific data sets to accompany each of the hands-on tutorials and encouraged learners to use these data sets in the classroom. Some data sets were used for more than one tutorial to provide continuity through the analysis workflow.

**Target learners and admissions.** We targeted our applicant pool at the learners who would benefit most from the training and who we expected would share their developed expertise with others to maximize the reach of the workshop's training. We accepted applicants who were novices in their analysis skillset and who did not have apparent access to other resources to support their skill development. We also aimed to promote diversity in scientific discipline (e.g., human, agricultural, environmental microbiomes), learner gender and background, research institution (e.g., undergraduate-serving, research university, agencies), geography (with special advertising directed to learners in the Midwest), and academic level (Fig. 2; see also Fig. 3). We also strove to provide opportunities to international learners and learners from underrepresented backgrounds. To advertise the course, we used social media (Twitter), our website, and professional networks. We also attempted to reach broader audiences by advertising with international scientific networks, especially Ciencia Puerto Rico in 2014 to 2016.

In each workshop, we could accommodate 23 to 26 learners in the classroom, and applications were oversubscribed every year (Table 1). As admissions became increasingly competitive, we began to require (rather than encourage) applicants to generate a microbiome data set prior to the workshop. We found that students who had struggled in an analysis attempt were highly incentivized to maximize their time at the workshop. Also, they could work on their data during office hours and ask specific questions of the instructors and teaching assistants (TAs).
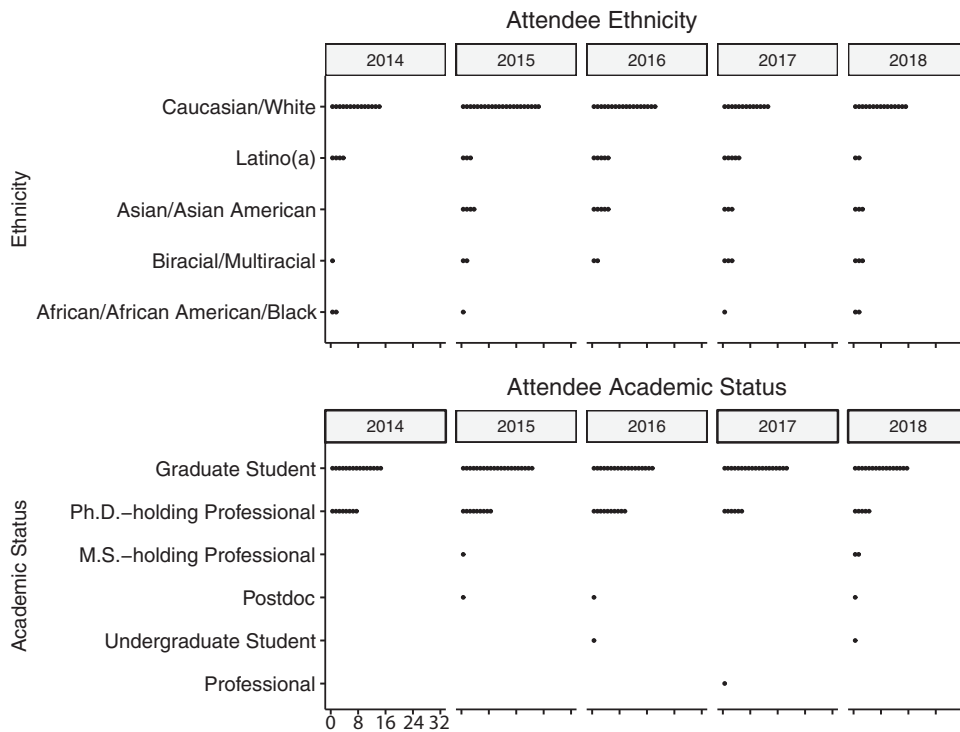
## Attendee Ethnicity



## Attendee Academic Status



**FIG 3** Distributions of EDAMAME learner ethnicity and academic status, 2014 to 2018. Data were collected and summarized from pre- and postworkshop assessments. The numbers at bottom left in the figure represent the number of participants, and it is the standard scale for all panels.

**Instructional team.** A large instructional team was necessary to support EDAMAME's learning goals. There were one to three lead instructors per year (Table 1). The instructors led the courses, oversaw admissions, provided lectures and course content, determined guest lectures, and mentored TAs in tutorial development. In the final 2 years of the workshop, there was also a course coordinator who oversaw conference logistics; fielded learner and applicant questions; and coordinated transportation for learners, guest lecturers, and instructors.

The hands-on nature of the workshop necessitated the presence of several dedicated TAs. Multiple instructors and supportive TAs in the classroom allowed us to be immediately responsive to the needs of the learners. TAs led tutorials based on interest and expertise. Having multiple TAs broadened instructional expertise and allowed unscheduled time for each TA to rest when they were not supporting instruction. Most often, new learners struggled with basic syntax and with interpreting error messages. Novice TAs (e.g., early graduate students) helped learners troubleshoot common errors,

**TABLE 1** Summary of EDAMAME dates and numbers of instructional staff, applicants, and learners from 2014 to 2018

| Yr | Dates | No. of days | No. of TAs | No. of instructors | No. of applicants | No. of workshop learners[a] |
|---|---|---|---|---|---|---|
| 2014 | 22 June to 29 June | 8 | 1 | 3 | 50 | 23 |
| 2015 | 21 June to 1 July | 11 | 6[b] | 1 | 93 | 32[c] |
| 2016 | 10 July to 20 July | 11 | 6 | 3 | 62 | 25 |
| 2017 | 6 August to 12 August | 7 | 7 | 3 | 63 | 26 |
| 2018 | 24 June to 30 June | 7 | 10 | 2 | 103 | 26 |

[a]The data representing workshop learners are from pre- and postsurvey responses. Additional local learners participated *ad hoc* and may not have completed surveys.

[b]There were two guest TAs in 2015 who participated in only one tutorial each, with the remaining 4 TAs available throughout the workshop.

[c]Participants in 2016 included 3 remote learners who participated in selected tutorials.

**TABLE 2** Guest lecturers and instructors for EDAMAME

| Yr | Guests |
|---|---|
| 2014 | C. Titus Brown (then at Michigan State University, now at the University of California, Davis); Jack Gilbert (University of Chicago); Pat Schloss (University of Michigan); Jim Tiedje (Michigan State); Sebastian Boisvert (Argonne National Laboratory); Stuart Jones (University of Notre Dame); Jay Lennon (Indiana University); Adina Howe (Michigan State University); Kathryn Docherty (Western Michigan University); Ariane Peralta (East Carolina University) |
| 2015 | Vince Young (University of Michigan); Pat Schloss lab members (University of Michigan); Ariane Peralta (East Carolina University); Jay Lennon (Indiana University); Stuart Jones (University of Notre Dame); Jim Tiedje (Michigan State); Jim Cole (Michigan State); Qiong Wang (Michigan State); Matt Scholz (Michigan State); Sarah Evans (Kellogg Biological Station); Vincent Denef (University of Michigan) |
| 2016 | Sarah Evans (Kellogg Biological Station); Pat Schloss (University of Michigan); Stuart Jones (University of Notre Dame); Jim Tiedje (Michigan State); Jim Cole (Michigan State); Rich Lenski (Michigan State University); Pat Bills (Michigan State University) |
| 2017 | Stuart Jones (University of Notre Dame); Pat Schloss (University of Michigan); Jim Tiedje (Michigan State University); Heather Allen (USDA, Ames, IA) |
| 2018 | Patrick Schloss (University of Michigan); Stuart Jones (University of Notre Dame); Tomas Vetrovsky (Czech Academy of Sciences); Thea Whitman (University of Wisconsin); Jim Tiedje (Michigan State University) |

while the more senior TAs and instructors assisted with more complicated problems (e.g., software and operating system incompatibilities, experimental design power for data analysis). In addition to instruction, TAs supported the logistical aspects of the course, such as providing local transportation for learners, purchasing supplies, and assisting learners with unexpected personal needs (e.g., trip to the medical center, forgotten toothbrush). TAs included volunteers (graduate students and postdocs) and graduate assistants who were partially supported by EDAMAME external funding. Participation in the workshop also offered TAs teaching opportunities that served diverse audiences.

There were also numerous invited guest instructors who offered tutorials, technical lectures, and research talks (Table 2). These guest instructors were selected to represent diverse career stages and levels of expertise across microbiome science (Table 2). Guest instructors varied according to guest availability, learner interests, and workshop duration, but some guest instructors generously provided content every year. Stuart Jones (University of Notre Dame) taught statistical analysis in R; Patrick Schloss and members of his lab (University of Michigan) taught amplicon analysis with mothur; Jim Tiedje (Michigan State University) provided a lecture and discussion on the future of microbial ecology. Instructors interacted with the learners during dinner and social time, and this provided an opportunity for learner networking and discussions.

**Learning environment and daily schedule.** EDAMAME was held at the Kellogg Biological Station (KBS), which offered a remote location and an immersive experience for learners and instructors. KBS was also chosen for economy—the room and board rates at KBS were affordable to many (e.g., ~$370 per week in 2018). Teaching assistants and volunteers provided transportation from the Kalamazoo and Lansing airports to KBS. KBS also provided conference services, dining, Wi-Fi, and bonfires. Finally, the natural setting and outdoor activities at KBS provided a respite from the time spent in front of the computer.

The durations of the workshops ranged from 7 to 11 days (Table 1; see also Fig. 4), including travel days. The morning schedule included an overview lecture followed by hands-on tutorials and group learning activities. After lunch, we had an afternoon lecture and additional tutorials. We held optional office hours with "choose your own adventure" tutorials and/or lectures on learner-chosen topics during the afternoon break. For example, in 2018 we discussed exact sequence variant analysis because it was a new approach being used in the field for defining operational taxonomic units. Learners could also ask specific questions about their own data during office hours. After dinner, we held an evening guest lecture in microbiome research. Evenings provided free time for networking and relaxation.

| TIME | DAY 1 | DAY 2 | DAY 3 | DAY 4 | DAY 5 | DAY 6 | DAY 7 |
|---|---|---|---|---|---|---|---|
| 8:00 AM | | Breakfast | Breakfast | Breakfast | Breakfast | Breakfast | |
| 8:30 | | | | | | | |
| 9:00 | | Welcome lecture and assessment | BLAST | | Ecological Analysis | Sequencing Variants and Workflows | Invited Keynote |
| 9:30 | | | | | | | |
| 10:00 | | | | | | | |
| 10:30 | | | | | | | |
| 11:00 | | Computing Workflows | Fetching Data | | | | Course Closeup and Assessment |
| 11:30 | | | | | | | |
| 12:00 PM | | Lunch | Lunch | Amplicon Sequencing | Lunch | Lunch | |
| 12:30 | | | | | | | |
| 1:00 | | | | | | | |
| 1:30 | | Shell Programming | R Programming | | Statistical Analysis | Metagenomes | |
| 2:00 | | | | | | | |
| 2:30 | | | | | | | |
| 3:00 | | Moving Data | Data visualization | | | | |
| 3:30 | | | | | | | |
| 4:00 | | | | | | | |
| 4:30 | | Break and/or Office Hours | Break and/or Office Hours | | Break and/or Office Hours | Break and/or Office Hours | |
| 5:00 | | | | Break | | | |
| 5:30 | | | | | | | |
| 6:00 | | | | | | | |
| 6:30 | | Dinner | Dinner | Dinner | Dinner | Dinner | |
| 7:00 | | | | | | | |
| 7:30 | | | | | | | |
| 8:00 | Meet and Greet | Invited Talk or Social Gathering | Invited Talk or Social Gathering | Invited Talk or Social Gathering | Invited Talk or Social Gathering | Invited Talk or Social Gathering | |
| 8:30 | | | | | | | |
| 9:00 | | | | | | | |

**FIG 4** An example of an ideal 1-week schedule for EDAMAME. For a full schedule with links to 2018 EDAMAME tutorials and talks, see https://github.com/edamame-course/2018-Tutorials/wiki/Schedule-EDAMAME-2018. Content for all 5 years of the workshop are also available at http://www.edamamecourse.org/materials/.

**EDAMAME educational strategy and assessment.** EDAMAME's educational strategy addressed two training needs. First, we offered general training in the fundamentals of introductory computing (e.g., command line, scripting, cloud computing, bioinformatic workflows [8]). This equipped participants with the basic skills needed to independently execute their analyses. We also offered specific training to overcome hurdles particular to microbial metagenomic data analysis and advised on best practices for microbiome analysis. To iteratively assess these strategies, we used a combination of summative and formative assessments to determine participant learning gains.

For the summative assessments, we worked with educational consultants to develop online, anonymous surveys and to perform pre- and postworkshop assessments. These assessments evaluated student-reported learning gains and confidence in areas aligned with our learning objectives. Each learner created a password to preserve anonymity while allowing for linking the pre- and postsurvey responses. To maximize response rates, we provided dedicated time in the classroom to complete the surveys. The preassessment survey was completed on the first full day, and the postassessment survey was completed on the final day of the workshop. We updated the survey annually to reflect any new or changed learning objectives but maintained the structure to facilitate interannual comparisons. Results of the annual surveys guided the continued development of course materials and topics covered. In the early years of the workshop, we had consultants perform in-classroom observations and provide feedback to the instructors. Ultimately, we compiled the 5 years of pre- and postsurvey data and performed a longitudinal analysis.

In the pre- and postsurveys, learners were asked to indicate the extent to which they understood specific learning outcomes or skills covered in the course, with ratings (e.g., Strongly Disagree, Disagree, Agree, and Strongly Agree [Fig. 5]).

We also used "real-time" assessment during the workshop by replicating formative assessment strategies found to be effective in Software Carpentry workshops (5–7). Formative assessment is an approach where teachers use informal practices to assess student learning during the workshop to evaluate understanding and modify teaching if needed. We used red and green sticky notes and "minute cards" to get this real-time feedback from students. Each participant was given a green ("I'm doing okay") and a

| |
|---|
| I know how to process Illumina data |
| I understand what per_library_stats.py does |
| I know how to run R |
| I know the main differences in analyses offered by QIIME and mothur |
| I am familiar with .biom formatted files |
| I can name at least two different microbial metagenomic databases |
| I know what an R package is |
| I understand the structure of an OTU table |
| I know what a kmer is |
| I know the difference between alpha and beta diversity |
| I know how to visualize microbial metagenomic data |
| I know how to use metadata to guide community analyses |
| I know how to assemble shotgun metagenomic data |

**FIG 5** Representative survey statements used for assessment in the "Computational Understanding" scale. OTU, operational taxonomic unit. References: QIIME (3), mothur (4), the .biom format (9), and R (10).

red ("I have a question") sticky note to stick onto their open laptop during tutorials. This visual cue allowed instructors to quickly survey the classroom and determine learners' comfort level and to attend to any student who was struggling during tutorials. Furthermore, it allowed students to continue working through tutorials or trouble-shooting without the need of raising their hand. We also employed minute cards. After each tutorial, students wrote what went well on the green sticky note and what could be improved on the red sticky note. Instructors and TAs read through notes during breaks to quickly identify gaps in understanding. This allowed us to identify gaps and make adjustments (e.g., in speed) in the subsequent instruction period.

**Building community resources and peer networks.** We were dedicated to promoting a welcoming and supportive learning environment. We presented a Code of Conduct in the welcome lecture, which outlines expectations for student and faculty behavior during the workshop and reporting procedures if those behavior expectations are not met, so that it was clear that any questionable conduct was grounds for dismissal. We created an online shared document Etherpad for collaborative note taking to maximize engagement and inclusivity. All materials were regularly updated and available online through our course website. We did our best to accommodate learners with families, providing private housing to families and learners with special requirements.

We aimed to build a peer learning community and to provide resources to support learners beyond the workshop. We offered an informal meet-and-greet on the arrival travel day and get-to-know-you short talks as lighting presentations after the first full day. These interactions allowed learners to identify peers with common research interests early in the workshop. We created a workshop website and a public repository on GitHub so that learners (and outside parties) could access EDAMAME learning materials. Linked content included lectures, hands-on tutorials, and reference lists. These materials have been shared openly, with most content licensed as Creative Commons generic (e.g., CC-BY), so that all course registrants and anyone else interested could have access. We also shared group email lists and encouraged social media outreach via Twitter and blogging. An EDAMAME meet-up was also held at the International Society for Microbial Ecology 2016 meeting in Montreal, Canada.

**Pre- and postsurvey comparisons and qualitative interviews.** Ninety-seven percent of EDAMAME learners from 2014 to 2018 rated the workshop overall in the top evaluative categories, "good" to "very good" (Fig. 6). A comparison of pre- and postassessment learner-reported learning gains and/or confidence with the major learning objectives of EDAMAME shows gains in all subcategories of learning reported
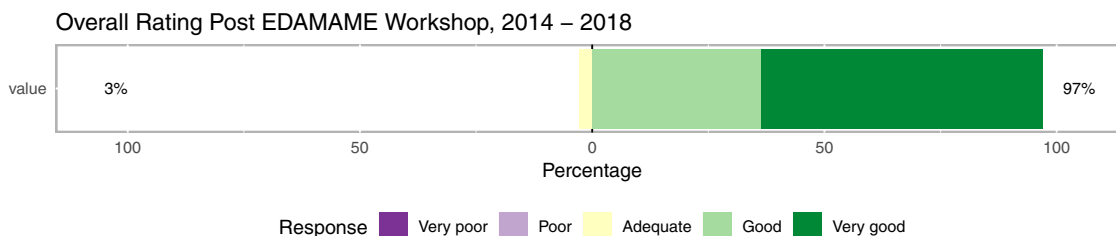
## Overall Rating Post EDAMAME Workshop, 2014 – 2018



**FIG 6** Overall EDAMAME assessment, 2014 to 2018. Data were collected and analyzed from pre- and postworkshop assessments.

(Fig. 7). The largest gains between the pre- and postassessments were seen with Computational Understanding (Fig. 7B) and Perception in Ability (Fig. 7C).

We also asked short-answer questions at the end of the survey, in which learners were asked to design an experiment and report how they would process and analyze microbial community high-throughput sequencing data. We observed increased sophistication in the responses to the short-answer questions from the pre- to postsurvey



**FIG 7** Summarized comparison of self-reported learning gains between pre- and postworkshop assessments, aggregated over 2014 to 2018. (A) Comfort with computational tasks. (B) Computational understanding. (C) Perception of computing ability. (D) Coding ability.

**TABLE 3** Representative comments from interviews[a]

| Comment category | Comment content |
|---|---|
| Positive overall | "EDAMAME was an inspiring introduction into microbiology. I thought the kind of analyses you could do with microbiology was really interesting. I really got pulled in on the data science part." (2014); "It was definitely one of the most effective workshops I've been to." (2016); "Very comprehensive, reached a lot of people from different backgrounds who were interested in analyzing microbial communities. I thought it provided a good survey of the tools that were available and it brought in some experts." (2014) |
| Content overwhelming | "I loved it, I had a blast. It was exhausting. It was a lot of fun, I learned a lot. I kind of felt overwhelmed." (2016); "I appreciated the workshop for its usefulness, it's a lot to take in. We need time to process. It's nice to have a bit of a breather. For someone who was new to the field like me. I needed a bit of time to digest." (2014); "It was pretty intense for me. I had never done any kind of code work before. This was really my first introduction..." (2015) |
| Career impact | "I can say that the course inspired me and put me on my path and inspired me to think about different ways to do analysis. They talked a lot about the different tools that were available." (2014); "It was a great workshop. It really helped me in my career path. It's opened a door for me to get into bioinformatics." (2015); "It really propelled my graduate school career and has pushed me... I took away the basic tools and I've been able to grow from that... I know how to make a pipeline. I know the basic structure and they gave that to me." (2015); "I'm one of the few people at [my workplace] who can analyze sequence data." (2016) |

[a]The sample was small at nine attendees, but each interviewee spent somewhere between 25 and 40 min discussing their experience at the workshop and its impact on their professional life and walking through the agenda for their year's workshop to give detailed feedback. While it is a small sample, each person contributed a lot of information. There were two respondents for 2014, four for 2015, and three for 2016. Each quote is labeled with the year in which the respondent participated in the workshop.

periods, with some learners leaving the questions blank in the presurvey and then providing thorough answers in the postsurvey (data not shown, but anonymized annual assessment reports are available upon request). This suggests large gains in particular for those learners who were new to high-throughput sequence analysis.

Qualitative interviews with 9 learners who attended EDAMAME from 2014 to 2016 (each spending 25 to 40 min with the interviewer; Table 3) suggested that the members of this group of learners were largely satisfied with the workshop and appreciated the attentiveness of the TAs and instructors as well as the red/green sticky note mechanism for soliciting help in real time. However, some of these learners also felt that there was too much material covered in the workshop and reported that they struggled to keep up with the pace of the course ("Content overwhelm"). Finally, we had many interviewed learners state that the workshop and materials covered made a positive impact on their career and research.

## DISCUSSION

**Lessons learned.** We offer suggestions from our experiences for running an effective microbiome analysis workshop (Fig. 8). EDAMAME's content changed from 2014 to

1. Regularly evaluate and change content to meet changing learner needs.
2. Maintain a high instructor to learner ratio.
3. Provide consistent workshop timing and fill the "middle-ground" duration needs of learners.
4. Understand the pros and cons of cloud computing for a workshop, and plan use of these resources well in advance.
5. Reach the broadest applicant pool of learners who have the potential to have the most gains from the training.
6. Consider the trade-off in workshop value (including instructor time) and maintaining economical costs to learners.
7. Plan well in advance to achieve best outcomes for applicants who require a US visa and international travel plans to attend the workshop.
8. Almost all learning engagement needs to happen on-site; efforts to engage learners pre-workshop were ineffective.
9. Scheduled classes should teach to the majority of learners to accomplish all our learning objectives. Office hours can help struggling learners catch up.
10. A welcoming and inclusive environment creates a positive workshop experience and is essential for effective learning.

**FIG 8** Lessons learned over 5 years of the EDAMAME workshop.

2018 to meet changing learner needs. These changes were guided in part by the applicants' responses to questions about their data sets and their expectations for the workshop. For example, amplicon analysis (e.g., 16S rRNA gene sequencing) was favored in the early years whereas untargeted metagenome analysis was favored in later years. Similarly, proportionally fewer students in 2018 were novices with respect to the command line or R, but the majority of the class appreciated the refresher. Some of the learners with self-taught experience embraced the opportunity to relearn the "correct" approaches and to gain missing foundational knowledge. Several tutorials were popular every year. For example, there was a consistent demand for ecological statistics and "supporting" skills such as GitHub/version control and cloud computing.

High instructor-to-learner ratios were essential for the success of the hands-on EDAMAME workshop. In the years in which we had the lowest instructor-to-learner ratios (e.g., in 2014 and 2015; Table 1), the TAs and instructors anecdotally reported exhaustion whereas the learners craved more attention. In addition to assistance from the formal instructors, learners could assist one another. To facilitate peer learning, we arranged the classroom in tables with groups of two or four. We also encouraged learners to support one another with troubleshooting while waiting for an instructor to become available to provide assistance.

Regardless of the length of the course, several learners indicated in their postassessments that more time at the workshop was needed each year. However, learners who were faculty or staff researchers shared (in informal conversations) that they would have been unable to commit to a longer workshop due to other professional responsibilities. We noted that there were other offerings for multiweek workshops (e.g., STAMPS), as well as several 1-day or 2-day workshops at professional society meetings and pipeline-specific training (e.g., in mothur and QIIME).

The timing of the workshop had several challenges. EDAMAME was held in the summer, and we tried to avoid scheduling it for the same week as major microbiology conferences, such as the American Society for Microbiology Microbe meeting and the International Symposium on Microbial Ecology (ISME) and Ecological Society of America meetings. Because microbiome analysis spans multiple disciplines, it was hard to avoid all of the large conferences that microbiome researchers may attend. We also had to change the timing of the workshop every year to accommodate the KBS event schedule. As EDAMAME grew in popularity, some learners applied for fellowships or travel awards to support their training, but the annual changes in timing made it difficult for students to plan. Moving the workshop to a dedicated conference site (e.g., a hotel) might help with consistent timing, but it would also increase the cost to learners.

We found that using cloud computing streamlined the course content and democratized access. We used the Amazon Elastic Compute Cloud (EC2), which was cost-effective and available to students who do not have access to high-performance computers at their home institutions. In the early years of the workshops, we guided learners through software installation on the EC2, but in the later years, we installed the software on the EC2 for the learners so that they could focus on moving data to and from the EC2. Using the EC2 presented a challenge for learners who were affiliated with government agencies or research laboratories (e.g., the U.S. Environmental Protection Agency, U.S. Geological Survey) because of their need for additional security and management approval prior to installing new software or moving data. While we did not have a perfect solution for these learners, we began to anticipate their needs and prompted them to apply for the required permissions in advance. Another hurdle with using the EC2 was the changing ways in which Amazon provided student or educational computing resources over the years. Amazon provided individual credits to learners in some years and in others required the instructors to apply for an educational grant. Cloud computing logistics needed to be anticipated about 9 months in advance, but in years where individual email addresses were needed, it was impossible to prepare until after the admissions were finalized, which typically occurred 4 to 6 months in advance of the workshop. We also noted an issue for some international

learners who did not have credit cards compatible with Amazon requirements to enroll for an EC2 account, and for these learners we had to share our own accounts or create accounts for them.

While our applicant pool and learner demographics reflected balance in sex, discipline, and academic level, EDAMAME fell short of its racial diversity goals. We could have benefitted from improvement in advertising the course to reach a broader pool to attract more applicants of color. We largely advertised on social media and through word of mouth. We attempted to specifically advertise to key target learner groups, such as those underrepresented in the sciences who may be expected to have less access to the training. On a positive note, we have evidence that EDAMAME was reaching socioeconomic diversity goals, as two interview respondents were clear that they would not have had the same opportunity for training and advancement given their lower-income backgrounds if it had not been for EDAMAME.

A final lesson to share is the balance between course value and learner costs. In its first years, EDAMAME was funded piecemeal by generous sponsors. We experimented with a mixed enrollment model of offering EDAMAME for university credit to local students and for a fee to outside students, but many of the local students could not afford the summer tuition required for the credit hours. Subsequently, EDAMAME was funded by external federal grants (the NIH during 2015 to 2018 and the USDA during 2017 to 2018), and we stopped offering it for credit for logistical simplicity. We began to charge modest workshop fees ($325) to support items that could not be covered by the grant (e.g., coffee, snacks). As soon as we began to charge workshop fees, the majority of applicants began to request financial aid. We realized that many of the learners, mostly graduate students and postdocs, were paying for the workshop personally, so we then worked to waive fees for eligible students in need and offer scholarships for students with international travel. The instructional team did not have enough funds to fully pay the TAs and instructors, who largely volunteered their time because they believed in the mission of the training. Guest instructors and lecturers generously volunteered their time as part of their broader impact for federal grants, and the workshop covered their travel expenses along with room and board at KBS. Thus, there is inevitable tension balancing instructor compensation and course affordability.

How much does it actually cost to run a workshop like EDAMAME? The first year, we ran the workshop for less than $14,000; students paid their own expenses of room and board, and no workshop fees were charged. This face amount does not include the substantial additional support that was provided via shared logistics with the ANGUS workshop, which was occurring at the same time at Kellogg Biological Station. It also does not include any support for personnel, which was the largest expense. Ideally, there would have been an annual budget for instructor and TA summer salaries, a logistics coordinator salary, and an hourly salary for undergraduate labor during the course. We also realized that unless we could procure funds to support personnel, the training might not be valued as highly by institutions and peers and might instead be perceived as a representing a cost diverting funds from other scholarly activities. The second biggest expenses were those of financial aid to offset costs of room and board and workshop fees to learners who needed it, which we provided in 2017 and 2018 to qualified learners, with USDA support. The third biggest expense was represented by the educational consultant employed to evaluate the course as a neutral third party and totaled $5,500 to $6,000 per evaluation. The remaining expenses were those represented by conference services at Kellogg Biological Station and lodging and travel expenses for the instructional team and guest speakers. In summary, there is a trade-off between the course cost, inclusive of the real value of instructor/TA time, and workshop affordability for the learners.

**Future directions.** While the data indicate that EDAMAME workshop was effective, only a limited number of learners can be accommodated per year, and there is high effort from the instructional team to support them. This is a low-throughput model of skill development. We are eager to reach a larger learner pool than what we could

accommodate in the classroom. In 2016, we experimented with live engagement of three to five remote learners (the number varied by tutorial) using free conference calling and screen sharing resources. The remote learners participated as a group at the same location. They engaged with the lectures and tutorials as fully as possible (but missed out on the guest lectures and other events). This added a mild distraction for the on-site learners, but the workshop proceeded relatively smoothly. The biggest hurdle was engaging with the remote learners during tutorials, as they had no classroom support. It is possible that a remote learning workshop could be successful, given an appropriate investment into conference technology, an on-site coordinator dedicated to its logistics, and an enhanced instructional team with traveling TAs dedicated to the remote classrooms.

The content of EDAMAME remains freely available online (https://github.com/edamame-course), but parts of the content are also being transitioned to local offerings. Many universities desire more offerings of online or digitized curriculum, and there is a question of how to balance the university's need to provide quality instruction for tuition with the open-science philosophy of providing free, democratic access to information. At Michigan State University, we are developing a graduate-level learning module on microbial metagenomics that includes amplicon and untargeted metagenome analysis pipelines. The 1-credit metagenomics module includes hands-on tutorials, is offered twice a week for 1 month, and is accompanied by prerecorded lectures. Postdoctoral trainees or faculty can enroll for a modest fee. Though based on EDAMAME materials, the modular content at Michigan State covers less content because there are prerequisite modules required for enrollment. Learners already have familiarity with the command line, with submitting jobs to the high-performance computing cluster, and with fundamentals of microbial genome analysis. EDAMAME materials have also been extended to international workshops, including a metagenomics 1-day crash course in Rio, Brazil, and a 1-week microbiome analysis workshop at Centro de Investigaciones Biológicas del Noroeste in La Paz, Mexico. In addition, more general tutorials (e.g., shell, GitHub, etc.) remain available from other data science training efforts, including Software and Data Carpentry, and short-format 2-day workshops on these skills are available through The Carpentries (http://carpentries.org).

Finally, we seek to maximize the impact of EDAMAME by offering this kind of training to those who need it most. We hope that the impact of our trainees training others is a lasting legacy of EDAMAME. We have found that our international learners have benefited immensely from this course, as they are challenged by limitations of access to computer resources or training. Going forward, we hope to continue to identify target audiences who could both benefit from our training and extend its impact broadly. Additionally, sequence analysis will continue to evolve with technologies, impacting the depth and breadth of scientific questions and experiments that are imaginable. We hope that our course content can continue to remove obstacles for scientists who wish to engage in these technologies.

## MATERIALS AND METHODS

This research was reviewed and exempted by Michigan State University's Human Research Protection Program under IRB i052533 (standard educational practices), as reviewed by the Michigan State University Biomedical and Health Sciences Institutional Review Board (BIRB) and Social Science/Behavioral/Education Institutional Review Board (SIRB).

Data analysis for the pre- and postsurvey assessment and associated reports were generated by outside research consultants. Final reports for the years 2016, 2017, and 2018 were written by Beth M. Duckles of Insightful, LLC, and for the years 2014 and 2015, reports were written by Julie Libarkin of STEM E.D., LLC. The code is available at https://github.com/ShadeLab/EDAMAMESurveys. Beth M. Duckles of Insightful also conducted qualitative interviews and provided final demographic summaries.

## REFERENCES

1. Quince C, Walker AW, Simpson JT, Loman NJ, Segata N. 2017. Shotgun metagenomics, from sampling to analysis. Nat Biotechnol https://doi.org/10.1038/nbt1217-1211b.

2. Knight R, Jansson J, Field D, Fierer N, Desai N, Fuhrman JA, Hugenholtz P, van der Lelie D, Meyer F, Stevens R, Bailey MJ, Gordon JI, Kowalchuk GA, Gilbert JA. 2012. Unlocking the potential of metagenomics through replicated experimental design. Nat Biotechnol 30:513–520. https://doi.org/10.1038/nbt.2235.

3. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK, Fierer N, Peña AG, Goodrich JK, Gordon JI, Huttley GA, Kelley ST, Knights D, Koenig JE, Ley RE, Lozupone CA, Mcdonald D, Muegge BD, Pirrung M, Reeder J, Sevinsky JR, Turnbaugh PJ, Walters WA, Widmann J, Yatsunenko T, Zaneveld J, Knight R. 2010. QIIME allows analysis of high-throughput community sequencing data Intensity. Nat Methods 7:335–336. https://doi.org/10.1038/nmeth.f.303.

4. Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, Lesniewski RA, Oakley BB, Parks DH, Robinson CJ, Sahl JW, Stres B, Thallinger GG, Van Horn DJ, Weber CF. 2009. Introducing mothur: open-source, platform-independent, community-supported software for describing and comparing microbial communities. Appl Environ Microbiol 75:7537–7541. https://doi.org/10.1128/AEM.01541-09.

5. Teal TK, Cranston KA, Lapp H, White E, Wilson G, Ram K, Pawlik A. 2015. Data carpentry: workshops to increase data literacy for researchers. Int J Digit Curation 10:135–143. https://doi.org/10.2218/ijdc.v10i1.351.

6. Wilson G. 2006. Software carpentry: getting scientists to write better code by making them more productive. Comput Sci Eng 8:66–69. https://doi.org/10.1109/MCSE.2006.122.

7. Wilson G, Bryan J, Cranston K, Kitzes J, Nederbrag L, Teal TK. 2017. Good enough practices in scientific computing. PLoS Comput Biol 13: e1005510. https://doi.org/10.1371/journal.pcbi.1005510.

8. Shade A, Teal T. 2015. Computing workflows for biologists: a roadmap. PLoS Biol 13:e1002303. https://doi.org/10.1371/journal.pbio.1002303.

9. McDonald D, Clemente JC, Kuczynski J, Rideout J, Stombaugh J, Wendel D, Wilke A, Huse S, Hufnagle J, Meyer F, Knight R, Caporaso J. 2012. The Biological Observation Matrix (BIOM) format or: how I learned to stop worrying and love the ome-ome. GigaSci 1. https://doi.org/10.1186/2047-217X-1-7.

10. R Core Team. 2013. R: a language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria.