

Reintroducing mothur: 10 years later

Patrick D. Schloss^{1†}

† To whom correspondence should be addressed: pschloss@umich.edu

¹ Department of Microbiology and Immunology, University of Michigan, Ann Arbor, MI 48109

Observation format

¹ **Abstract**

² 75 words

3 **Importance**

4 150 words

5 3000 words total

6 Few scientists set out on a nearly two decade long journey with a specific goal in mind. Often we
7 fail to start a scientific journey because it looks too hard. Perhaps we get bogged down in all of
8 the things that could go wrong. Perhaps we go astray from the path because we find something
9 else that appears more interesting. Every scientist picks their own path and takes their own forks in
10 the road. From the outside, it may appear to be a random walk. Nevertheless, these meandering
11 journeys in are common in science.

12 At the risk of navel gazing, looking back on our scientific journeys can be instructive to other
13 scientists who are overwhelmed at the prospect of looking forward at their careers [Lenski and
14 Chemistry guy]. By no means is my personal journey over, but since 2002 I have been on a
15 journey that I did not realize I was on. Now that the paper introducing the mothur software package
16 is ten years old and has become the most cited paper published by *Applied and Environmental*
17 *Microbiology*, it is worth stepping back and using the development of mothur as a story that likely
18 has parallels to many other research stories that have taken time to develop.

19 I fondly recall preparing a poster for the 2002 meeting of research groups supported by the
20 NSF-supported Microbial Observatories Program. I wanted to triumphantly show that I had
21 sequenced more than 600 16S rRNA gene sequences from a single 0.5-g sample of Alaskan
22 soil. This was greater sequencing depth than anyone else had achieved for a single sample. As
23 I was preparing the poster, I walked into the office of Jo Handelsman, my postdoctoral research
24 advisor, and laid out the outline for the poster. She asked if I could add one of those “curvy
25 things”, a rarefaction curve, to show where I was in sampling the community. Rarefaction curves
26 and attempts to estimate the taxonomic richness of soil had become popular because of the
27 simple, but impactful mini-review by Jennifer Hughes and her colleagues, which introduced the
28 field to operational taxonomic units (OTUs), rarefaction curves, and richness estimates [DOI:
29 10.1128/AEM.67.10.4399-4406.2001]. I do not recall whether that poster had a rarefaction curve
30 on it, but her question primed my career.

31 ***Introducing DOTUR and friends.*** When Jo asked me to generate a rarefaction curve for the
32 poster, the request was not trivial. How would I bin the sequences into OTUs? Hughes and her

colleagues did it manually for datasets that had fewer than 284 sequences. Although I could possibly do that for my 600 sequences, my goal was to generate 1,000 sequences from the sample and to repeat that sampling effort for other samples. I needed something that could be automated. Furthermore, the software that Hughes used, EstimateS, required a series of tedious data formatting steps. I had found my first problem. How would I assign sequences to OTUs and use that data to estimate the richness and diversity of a sample? The second problem would be how could I compare the sequences found in one sample to another sample? The solution to the first problem, DOTUR (Distance-based OTUs and Richness), took us two years to develop [10.1128/AEM.71.3.1501-1506.2005]. DOTUR did two things: given a matrix describing the genetic distance between pairs of sequences, it would cluster those sequences into OTUs for any distance threshold to define the OTUs and then it would use the frequency of each OTU to calculate a variety of alpha diversity metrics. The solutions to the second problem would come from our work to develop software including S-LIBSHUFF [2004], SONS (Shared OTUs and Similarity) [2006], and TreeClimber [2006]. Around the same time, Catherine Lozupone and Rob Knight were developing their UniFrac tools to compare communities with a phylogenetic rather than OTU-based approach [PMID: 16332807; PMID: 17220268]. With these tools, the field of microbial ecology had a quantitative toolbox for describing and comparing microbial communities.

It is important to remember that we knew there were many problems with 16S rRNA gene sequencing. We knew there were biases from extractions and amplification [XXXX]. We knew there were chimeras [XXXX]. Getting to the distance matrix required trimming and correcting sequence errors, aligning them, and finding the most appropriate way to calculate a pairwise distance. It was also rare to have experimental replication to perform statistical tests to compare treatment groups. We frequently used a dataset comparing Scottish soils from Alison McCaig and colleagues. This dataset consisted of two experimental groups, each replicated three times with 45 sequences in each replicate. Nevertheless, we had excuses and work arounds for these problems that served our needs. At the time, I felt that the biggest problems were how to cluster the sequences into OTUs and how to use those clusterings to test our hypotheses. Along the way we would demonstrate the utility of such tools to answer questions like where are we in the bacterial census? How many sequences would it take to see every OTU in that sample of Alaskan soil? How does the word

usage of *Goodnight, Moon* compare to that of *Portrait of a Lady*? More importantly, 1,900 papers used DOTUR to facilitate their own research questions. Had we waited to solve all of the problems that plague 16S rRNA gene sequencing, we would still be waiting.

As we developed these tools, I found a unique niche in microbiology. I believe that my undergraduate and graduate training as a biological engineer prepared me to think about research questions from a systems perspective, to think quantitatively, and to understand the value of using computer programs to help solve problems. As an undergraduate engineering student, I learned the Pascal programming language and promptly forgot much of it as an engineering graduate student. As a postdoc, I learned the Perl programming language to better understand how LIBSHUFF, a tool for comparing the structure of two communities, worked since it was written in Perl [DOI: 10.1128/AEM.70.9.5485-5492.2004]. After writing my own version of LIBSHUFF and seeing the speed of the version written in C++ by my collaborator, Bret Larget, I converted my Perl version of DOTUR into C++. At the time, the conversion from Perl to C++ seemed like an academic exercise to learn a new language. My Perl version only took a minute or so to process the final collection of 1,000 sequences and the C++ version took seconds. Was that really such a big difference? In hindsight, as we now process datasets with millions of sequences, the decision to learn to C++ was critical. The ability to pick up computer languages to solve problems was enabled by my prior training. It was also a skill that was virtually unheard of in microbiology. Today researchers without the ability to program in Python or R are at a significant disadvantage.

Introducing mothur. Shortly after DOTUR was published, I received an email from Mitch Sogin asking whether DOTUR could handle more than a million sequences. Without answering his question, I asked where he found a million sequences. Little did I know that his email would represent another pivot in the development of these tools. His group would be the first to use 454 sequencing technology to generate 16S rRNA gene sequences [PMID: 16880384]. Although mothur could assign those sequences to OTUs, it was slow and required a significant amount of RAM. As I left my postdoc to start my independent career across the state from Sogin's lab at the University of Massachusetts in Amherst, my plan was to rewrite DOTUR, SONS, S-LIBSHUFF, and TreeClimber for the new world of massively parallelized sequencing. The new tool was mothur.

Milling about at a poster session at an ASM General Meeting in New Orleans, I again ran into Mitch who asked what my plans were for new tools. I told him that I wanted to make a tool like ARB (a powerful database tool and phylogenetics package), but for microbial ecology analysis. His retort was, “You and what army?” To that point, I had written every line of code and been answering many emails from people asking for help. It would be difficult, but I needed to learn to let go and share the development process with someone else. He was right, I would need an army. That “army” ended up being Sarah Westcott who has worked on the mothur project largely from its inception. Today, mothur is over 200,000 lines of code and Sarah has touched or written nearly every line of code. Beyond writing and testing mothur’s code base, she has become a conduit for many learning the tools of microbial ecology by patiently answering questions via email and the package’s discussion forum. The community and I are lucky that Sarah has stayed with the project for more than a decade. To be honest, such dependency on a single person makes the project brittle. In hindsight, it would have been better to have developed mothur with more of an “army” or team so that there is overlap in people’s understanding of how mothur works. Although such an arrangement might work in a software engineering firm, it is not practical in an academic setting where funding is limited for developers of free, open source software packages. There are certainly projects that make this work, but they are rare.

Challenges of making open source count. Anyone can post code to GitHub with a permissive license and claim to be an open source software developer. Far more challenging is engaging the target community to make contributions to that code. Frankly, we have struggled to expand the number of people that make contributions to the mothur code base. One challenge we face is that if we looked to third parties to contribute code to mothur, they would need to know C++. Given the paucity of microbiologists with any programming skills, expecting that community to provide contributors that can write code in a syntax that prizes execution efficiency over developer efficiency was not likely. In contrast, the QIIME development team could be more distributed because their code base was primarily written in Python, which prizes developer efficiency over execution efficiency and exists as a series of wrappers to execute other developers’ code. These choices resulted in many tradeoffs that have impacted ease of installation, usability, execution speed, and flexibility. If we were offered a grant to rewrite mothur, we would likely rewrite it as an

R package that leaned heavily on the R language's C++ packages. Of course, such choices are always best in hindsight and when we started developing mothur, the ability to interface between scripting languages like R and Python and C++ code was not as well developed as it is today. For example, the modern version of the Rcpp package was first released in 2009 and its popularity was not immediate. Again, the development of mothur has been a product of the environment that it was created in. Although these decisions have largely had positive outcomes, there have been tradeoffs that caused us to sacrifice other goals.

Beyond contributing to the mothur code base, we sought out other ways to include the community as developers. The paper describing mothur included **N** co-authors, all but three (Schloss, Ryabin, and Westcott) responded to a call to provide a wiki page that described how they used an early version of mothur to analyze a data set. Our vision was that authors might use the mothur wiki to document reproducible workflows for papers using mothur but to also provide instructional materials for other seeking to adapt mothur for their uses. Again, this vision was a product of the environment. IPython notebooks (2011) and R markdown (2012) would not be developed until later. Unfortunately, once the incentive of co-authorship was removed, researchers stopped contributing their workflows to the wiki. Part of the difficulty of recruiting wiki contributors was a perception by some that the wiki was not a community resource. For example, I would frequently receive emails from people telling me that there was a typo on a specific page when the intention was that they could correct the typos without my input. We have been more successful in soliciting input and contributions from the user community through the mothur discussion forum and GitHub issue tracker. As mothur has matured, we have been dependent on the user community to use these resources to tell us what features they would like to see included in mothur. The user community also tells us where our documentation is confusing. Often we can count on people not directly affiliated with mothur to provide instruction and their own experience to other users. We are constantly trying to recruit our "army" and are happy to take any contributions we can. Whether the contributions are to the code base, discussion forum, or suggestions for new tools, these contributions have been invaluable to the growth and popularity of mothur.

Failed experiments. If we never failed, we would not be trying hard enough. Over the past decade we have tried a number of experiments to improve the usability and utility of mothur. One of our

first experiments was to use mothur to generate standard vector graphic (SVG)-formatted files of heatmaps and venn diagrams depicting the overlap between microbial communities. I quickly realized that I would never put a mothur-generated figure into a manuscript I wrote. Such visuals require far too much customization to be publication-quality. Although QIIME has incorporated visualization tools through the Emperor package, the challenge of users taking default values has downsides as ordinations with black background or publishing 3-D ordinations in a 2-D medium litter the literature. Instead, we have encouraged users to use R packages to visualize mothur-generated results using the minimalR instructional materials. A second experiment we attempted was to create a graphical user interface (GUI) for running mothur. Forcing users to interact with mothur through the command line has been a significant hurdle for many. Unfortunately, the development effort required to create and maintain a GUI is significant and there is limited funding for such efforts. The newest version of QIIME (version 2) has emphasized interaction with the tools through a GUI and it remains to be seen how this experiment will go. Another downside of using a GUI is that there is a risk that reproducibility will suffer if users do not have a mechanism to document their mouse clicks. This documentation is explicit in mothur as all commands and output is recorded in a logfile. Given the heightened recent focus on reproducibility we have extended significant effort in developing instructional materials teaching users how to organize, document, and execute reproducible pipelines that allow a user to go from raw sequence data to a compiled manuscript with figures through the Riffomonas project. Finally, we collaborated with programmers through Google Summer of Code to develop commands in mothur to implement the random forest and SVM machine learning algorithms. Similar to the challenges of developing attractive visuals, fitting the algorithms' hyperparameters, testing, and deploying the resulting models requires a significant amount of customization. Furthermore, this is an active area of research where methods are still being developed and improved. Thankfully, there are numerous R and Python packages that do a better job of developing these models. Again, we have put our efforts into developing instructional materials that mothur users can use to fit such models to their data. In each of our "failed" experiments, the real problems were straying from what mothur does well and failing to grasp what we really wanted the innovation to do. In hindsight, our solution to these failure has been to provide tutorials to as a conduit between mothur and their goals.

177 ***Competition is good and healthy.***

- 178 • A lot of people have been working in this space - e.g. clustering algorithms
- 179 • Approaches to developing software have differed (e.g. web-based (greengenes/RDP),
- 180 boutique (Quince), wrappers (QIIME))

181 What are we proud of? * Open source packages * Platform independent * Development of open
182 instructional materials * Data driven development of methods * Database independent - allows
183 people outside of our field to use mothur * Reproducibility * Helped create standards - when
184 something is this popular, you force people to use a certain suite of tools and methods

185 The future * 16S rRNA gene sequencing is not dead * Always improving sequencing methods
186 bring new problems * Larger datasets, need a replacement for naive open and closed reference
187 clustering * Lingering controversies - singletons, ASVs, contamination

188 Longterm funding * The development of mothur was enabled by initial funding from a grant to
189 Sogin from the Sloan Foundation to support his VAMPS (Visualization and Analysis of Microbial
190 Population Structures) initiative. * NSF * NIH * Bootstrapping from other projects

191 Red Queen

193 **References**

194 25 references

195 **Figure 1. Caption caption caption.** Footnotes footnotes footnotes.