# Adapterama II: Universal amplicon sequencing on Illumina platforms (TaggiMatrix)

Travis C. Glenn[1–4], Todd W. Pierson[1,§], Natalia J. Bayona-Vásquez[1,4], Troy J. Kieran[1], Sandra L. Hoffberg[2,*], Jesse C. Thomas IV[1, α], Daniel E. Lefever[5, ‡], John W. Finger Jr.[1,3, θ], Bei Gao[1, φ], Xiaoming Bian[1,ε], Swarnali Louha[4], Ramya T. Kolli[3, 6 ,¥], Kerin Bentley[2, #], Julie Rushmore[7,γ], Kelvin Wong[8,ω], Timothy I. Shaw[4,8,¶], Michael J. Rothrock Jr.[9], Anna M. McKee[10], Tai L. Guo[5], Rodney Mauricio[2], Marirosa Molina[8,Ω], Brian S. Cummings[3,6], Lawrence H. Lash[11], Kun Lu[1, ψ], Gregory S. Gilbert[12, 13], Stephen P. Hubbell[13, 14], & Brant C. Faircloth[15]

[1] Department of Environmental Health Science, University of Georgia, Athens, GA 30602, USA
[2] Department of Genetics, University of Georgia, Athens, GA 30602, USA
[3] Interdisciplinary Toxicology Program, University of Georgia, Athens, GA 30602, USA
[4] Institute of Bioinformatics, University of Georgia, Athens, GA 30602, USA
[5] Department of Veterinary Biosciences and Diagnostic Imaging, University of Georgia, Athens, GA 30602, USA
[6] Department of Pharmaceutical and Biomedical Sciences, University of Georgia, Athens, GA 30602, USA
[7] School of Ecology & College of Veterinary Medicine, University of Georgia, Athens, GA 30602, USA
[8] US Environmental Protection Agency, Athens, GA 30605, USA
[9] USDA-ARS U.S. National Poultry Research Center, Athens, GA 30605, USA
[10] U.S. Geological Survey, South Atlantic Water Science Center, Norcross, GA 30093, USA
[11] Department of Pharmacology, Wayne State University, Detroit, MI 48201, USA
[12] Environmental Studies Department, University of California, Santa Cruz, CA 95064, USA
[13] Smithsonian Tropical Research Institute, Balboa, Ancon, Republic of Panama
[14] Department of Ecology and Evolutionary Biology, University of California, Los Angeles, CA 90095, USA
[15] Department of Biological Sciences and Museum of Natural Science, Louisiana State University, Baton Rouge, LA 70803, USA
§ current address: Department of Ecology and Evolutionary Biology, University of Tennessee, Knoxville, TN 37996, USA
* current address: Department of Ecology, Evolution, and Environmental Biology, Columbia University, New York, NY 10027, USA
α current address: Division of STD Prevention, Centers for Disease Control and Prevention, Atlanta, GA 30329, USA
‡ current address: Integrative Systems Biology and Drug Discovery Institute, University of Pittsburgh, Pittsburgh, PA 15260, USA

39  θ current address: Department of Biological Sciences, Auburn University, Auburn, AL 36849,
40      USA
41  ϕ current address: Department of Medicine, University of California, San Diego, CA 92093, USA
42  ε current address: Complex Carbohydrate Research Center and Department of Microbiology,
43      University of Georgia, Athens, GA 30602, USA
44  ¥ current address: Epigenetics and Stem Cell Biology Laboratory, National Institute of
45      Environmental Health Sciences, Research Triangle Park, NC 27709, USA
46  # current address: LeafWorks Inc., 125 South Main Street #150, Sebastopol, CA 95472, USA
47  γ current address: Epicenter for Disease Dynamics, One Health Institute, School of Veterinary
48      Medicine, University of California, Davis, CA 95616, USA
49  ω current address: California Water Service, 1720 N First St, San Jose, CA 95112, USA
50  ¶ current address: Department of Computational Biology, St. Jude Children's Research Hospital,
51      Memphis, TN 38105, USA
52  Ω current address: National Exposure Research Laboratory, US Environmental Protection
53      Agency, Research Triangle Park, NC 27709, USA
54  Ψ current address: Department of Environmental Sciences and Engineering, University of North
55      Carolina, Chapel Hill, NC 27599, USA
56
57  Corresponding Author:
58  Travis Glenn
59  Dept. of EHS, Environmental Health Science Bldg., University of Georgia, Athens, GA 30602,
60  USA
61  Email address: travisg@uga.edu
62

## Abstract

64  Next-generation sequencing (NGS) of amplicons is used in a wide variety of contexts. Most
65  NGS amplicon sequencing remains overly expensive and inflexible, with library preparation
66  strategies relying upon the fusion of locus-specific primers to full-length adapter sequences with
67  a single identifying sequence or ligating adapters onto PCR products. In *Adapterama I*, we
68  presented universal stubs and primers to produce thousands of unique index combinations and a
69  modifiable system for incorporating them into Illumina libraries. Here, we describe multiple
70  ways to use the *Adapterama* system and other approaches for amplicon sequencing on Illumina
71  instruments. In the variant we use most frequently for large-scale projects, we fuse partial
72  adapter sequences (TruSeq or Nextera) onto the 5' end of locus-specific PCR primers with
73  variable-length tag sequences between the adapter and locus-specific sequences. These fusion
74  primers can be used combinatorially to amplify samples within a 96-well plate (eight forward
75  primers + 12 reverse primers yield 8 x 12 = 96 combinations), and the resulting amplicons can be
76  pooled. The initial PCR products then serve as template for a second round of PCR with dual-
77  indexed iTru or iNext primers (also used combinatorially) to make full-length libraries. The
78  resulting quadruple-indexed amplicons have diversity at most base positions and can be pooled

79    with any standard Illumina library for sequencing. The number of sequencing reads from the
80    amplicon pools can be adjusted, facilitating deep sequencing when required or reducing
81    sequencing costs per sample to an economically trivial amount when deep coverage is not
82    needed. We demonstrate the utility and versatility of our approaches with results from six
83    projects using different implementations of our protocols. Thus, we show that these methods
84    facilitate amplicon library construction for Illumina instruments at reduced cost with increased
85    flexibility. A simple web page to design fusion primers compatible with iTru primers is available
86    at: http://baddna.uga.edu/tools-taggi.html. A fast and easy to use program to demultiplex
87    amplicon pools with internal indexes is available at: https://github.com/lefeverde/Mr_Demuxy.
88

## Introduction

90    Next-generation DNA sequencing (NGS) has facilitated a wide variety of benefits in the life
91    sciences (Ansorg, 2009; Goodwin, McPherson & McCombie, 2016), and NGS instruments have
92    an ever-growing capacity to generate more reads per run. Substantial progress has been made in
93    developing new, lower-cost instruments, but much less progress has been made in reducing the
94    cost of sequencing runs (cf., Glenn, 2011 vs. Glenn, 2016). Thus, the large number of reads from
95    a typical NGS run comes with a relatively large buy-in cost but yields an extremely low cost per
96    read. Frustratingly, within every NGS platform, the lowest-cost sequencing kits have the highest
97    costs per read (Glenn, 2011; 2016). This creates a fundamental challenge: how do we efficiently
98    create and pool large numbers of samples so that we can divide the cost of high capacity NGS
99    sequencing runs among many samples, thereby reducing the cost per sample?
100    It is well known that identifying DNA sequences (commonly called indexes, tags, or
101    barcodes; we use the term indexes throughout) can be incorporated during sample preparation for
102    NGS (i.e., library construction) so that multiple samples can be pooled prior to NGS, thereby
103    allowing the sequencing costs to be divided among the samples (see Faircloth & Glenn, 2012 and
104    references therein). When sufficient unique identifying indexes are available, many samples,
105    including samples from multiple projects, can be pooled and sequenced on higher throughput
106    platforms which minimizes costs for all samples in the pool.
107    In many potential NGS applications, the number of desired reads per sample is limited, so
108    the cost of preparing samples for NGS sequencing becomes the largest component of the overall
109    cost of collecting sequence data. Thus, it is desirable to increase the number of low-cost library
110    preparation methods available. As the cost of library construction is reduced, projects requiring
111    fewer DNA sequences per sample become effective to conduct using NGS (e.g., if sample
112    preparation plus sequencing for NGS is < sample preparation plus sequencing on capillary
113    machines, then it is economical to switch).
114

115    *Previous NGS amplicon library preparation methods*
116    Amplicon library preparations for NGS have been integrating indexes for more than a decade
117    (e.g., Binladen et al., 2007; Craig et al., 2008). Early NGS strategies consisted of conducting
118    individual PCRs targeting different DNA regions from one sample and then pooling them

119    together. Then, full-length adapters would be ligated to each sample pool, providing sample-
120    specific identifiers. This approach has the advantage of being economical regarding amplicon
121    production, primer cost, and pooling of amplicons prior to adapter ligation, as well as being
122    ecumenical because the resulting amplicons can be ligated to adapters for any sequencing
123    platform. The downside of this first approach is that adapters must be ligated to the amplicons,
124    which is time-consuming, expensive, and error-prone, and which can introduce errors into the
125    resulting sequences. To avoid ligation of adapters to amplicons, most NGS amplicon sequencing
126    strategies have subsequently relied upon the fusion of locus-specific primers to full-length
127    adapter sequences and the addition of identical indexes to both 5' and 3' ends (e.g., Roche fusion
128    primers; Binladen et al., 2007; Bentley et al., 2009; Bybee et al., 2011; Cronn et al., 2012;
129    Shokralla et al., 2014). These strategies often use the whole sequencing run for amplicons only.
130    Illumina platforms have traditionally struggled to sequence amplicons because: 1) the platform
131    requires a diversity of bases at each base position (Mitra et al., 2015), which is easily achieved in
132    genomic libraries but not in amplicon libraries; and 2) read-lengths are limited, making the
133    complete sequencing of long amplicons challenging or impossible.
134         Several alternatives have been proposed to resolve the first issue (i.e., low base-
135    diversity). Users have typically added a genomic library (e.g., the PhiX control library supplied
136    by Illumina) to amplicon library pools to create the base-diversity needed, but this method
137    wastes sequencing reads on non-target (PhiX) library. Second, to solve the issue of limited read-
138    length, described above, custom sequencing primers can be used in place of the Read1 and/or
139    Read2 sequencing primer(s) (Caporaso et al., 2011). This method allows for longer effective
140    read-lengths by removing the read-length wasted by sequencing the primers used for
141    amplification (e.g., 16S primer sequences), but it can be very expensive to optimize custom
142    sequencing primers, costing hundreds of dollars for each attempt. Another alternative is to use
143    the amplicons as template for shotgun library preparations, most often using Nextera library
144    preparation kits (Illumina 2018a). A fourth method is to add heterogeneity spacers to the indexes
145    in the form of one, two, three (etc.) bases before the index sequence (e.g., Cruaud et al., 2017),
146    but because amplicons can contain repeats longer than the heterogeneity spacers, it is still
147    possible to have regions of no diversity. Thus, all of the proposed solutions have specific
148    limitations, and none are particularly economical for sequencing standard PCR products from a
149    wide range of samples, as is typical in molecular ecology projects.
150
151    *NGS amplicon needs*
152         In general, NGS has been widely adopted to sequence complex amplicon pools where
153    cloning would have been used previously (e.g., 16S from bacterial communities or viruses within
154    individuals). Such amplicon pools may have extensive or no length variation. Amplicons for
155    single loci from haploid or diploid organisms (with no length variation between alleles) are
156    typically still sequenced via capillary electrophoresis at a cost of about $5 USD per read. In
157    contrast to the high cost of individual sequencing reads via capillary instruments, >50,000
158    paired-end reads can be obtained for $5 USD on the Illumina MiSeq. Unfortunately, MiSeq runs

159    come in units of ~$2,000 USD for reads that total a length similar to that of capillary sequencing
160    (Glenn, 2016; paired-end (PE) 300 reads). Thus, it would be desirable to have processes that
161    allow users to: 1) pool samples from multiple projects on a single MiSeq run and divide costs
162    proportionately, and 2) prepare templates (i.e., construct libraries) at costs less than or similar to
163    those of traditional capillary sequencing.
164         Characteristics of an ideal system include: 1) use of universal Illumina sequencing
165    primers; 2) minimizing total sample costs, ideally to be below standard capillary/Sanger
166    sequencing; 3) minimizing time and equipment needed for library preparations; 4) minimizing
167    buy-in (start-up) costs; 5) eliminating error-prone steps, such as adapter ligation, 6) maximizing
168    the number of samples (e.g., ≥ thousands) that can be identified in a pool of samples run
169    simultaneously, 7) maximizing the range of amplicons that can be added to other pools (e.g.,
170    from <1% to >90%), and 8) creating a very large universe of sample identifiers (e.g., ≥ millions)
171    so that identifiers would not need to be shared among samples, studies, or researchers, even
172    when coming through large sequencing centers.
173         Single-locus amplicon sequencing represents one extreme example of the needs identified
174    above. In some scenarios, researchers may only be sequencing a single short, homogeneous
175    amplicon where ≥ 20x coverage is excessive. The cost of sequencing reagents for only 20 reads
176    of 600 bases on an Illumina MiSeq using version 3 chemistry, which generates ~20 million
177    reads, is <$0.01 USD (i.e., 1 millionth of the run). It is impractical to amass 1 million amplicon
178    samples for a single run. However, a small volume of dozens or hundreds of samples can be
179    easily added into a MiSeq run with other samples/pools that need the remaining of reads.  By
180    paying the proportional sequencing costs for such projects, the cost of constructing libraries and
181    conducting quality control on the libraries becomes the largest component of the total cost of
182    collecting NGS data. Having the ability to combine libraries of many different kinds of samples,
183    each with their own identification indexes, is critical to the feasibility of this strategy. We have
184    developed, and describe below, a system to meet most of the design characteristics enumerated
185    above.
186         In this paper, we focus on library preparation methods for amplicons. We introduce
187    TaggiMatrix, which is an amplicon library preparation protocol that is built upon methods
188    developed in *Adapterama I* (Glenn et al., 2019). This general method can be optimized for
189    various criteria, including the minimization of library preparation cost and reduction of PCR
190    bias. Briefly, by tagging both the forward and reverse locus-specific primers with different,
191    variable-length index sequences, and also by including indexes in the iTru or iNext primers, we
192    create quadruple-indexed libraries with high base-diversity, enabling the use of highly
193    combinatorial strategies to index, pool, and sequence many samples on Illumina instruments.
194

## Materials & Methods

196    *Methodological objectives*

197  Our goal was to develop a protocol that would help to overcome the challenges of amplicon
198  library preparation and fulfill the characteristics of an ideal system enumerated above. We extend
199  the work of Faircloth & Glenn (2012) and Glenn et al. (2019) to achieve these goals.
200
201  ***Methodological approach***
202  Illumina libraries require four sequences (P5 + Read1 and P7 + Read2; Fig. 1), and can
203  accommodate internal index sequences on each end, (i.e., P5 + i5 index + Read1 and P7 + i7
204  index + Read2; Fig. 1; Illumina Sequencing Dual-Indexed Libraries on the HiSeq System User
205  Guide; Glenn et al., 2019). The Read1 and Read2 sequences can be of two types—TruSeq or
206  Nextera—. Just as in *Adapterama I* (Glenn et al., 2019), we have designed systems for both.
207  Our overall approach is to make amplicons with fusion primers (Fig. 2) that can use iTru
208  or iNext primers described in *Adapterama I* (Glenn et al., 2019) to make full-length Illumina
209  libraries (Fig. 3a; Figs. S1 and S2). The resulting libraries always contain dual-indexes in the
210  standard indexing positions and may optionally contain additional internal indexes (Figs. 1–3;
211  Table 1; Illumina, 2018b). These indexes are recovered through the four standard separate
212  sequencing reactions generated by Illumina instruments when doing paired-end sequencing (Fig.
213  3b).
214  Although iTru and iNext primers facilitate quick and low-cost additions of dual-indexed
215  adapters, this still requires a separate PCR reaction (but, see Discussion). Thus, when hundreds
216  of amplicons are to be sequenced, it becomes economical to use additional internal indexes
217  (Table 1) so that amplicons can be pooled prior to the use of iTru or iNext primers (Figs. 1 and
218  2). This approach should work with a wide variety of primers (e.g., Table 2). Such combinatorial
219  indexing is designed to work in 96-well plate arrays but can be modified for other systems.
220  Typically, eight indexed fusion forward primers (A–H) and 12 indexed fusion reverse primers
221  (1–12) are designed and synthetized (File S1). Then, each DNA sample in each well of the 96-
222  well plate can be amplified with a different forward and reverse primer combination (File S1,
223  PCR_Set_up). These PCR products can be pooled and amplified using a similar combinatorial
224  scheme with tagged universal iTru/iNext primers in the second PCR (Table 3), enabling the
225  large-scale multiplexing of samples in one Illumina run (Table 4). Finally, because Illumina
226  MiSeq platforms have documented issues in the quality of Read 2, particularly in GC-rich
227  regions (Quail et al., 2012), fusion primers can be designed to swap forward and reverse primers
228  with Read1 and Read2 fusions (e.g., R1Forward + R2Reverse, vs. R1Reverse + R2Forward;
229  "flipped" primers) to account for this issue (Fig. 2). It is also possible to do replicate
230  amplification with both sets of primers (regular and flipped), to significantly increase base
231  diversity in amplicon libraries.
232
233  ***TaggiMatrix applied case studies***
234  We tested iTru primers designed as described above in five different experiments covering a
235  wide range of experiments typically done in molecular ecology projects, and we tested iNext
236  primers designed as described above in a single project (Table 4). In each experiment, we used at

237    least two sets of primers: the first set (i.e., locus-specific fusion primers) generated primary
238    amplicons, and the second set (i.e., iTru or iNext) converted primary amplicons into full-length
239    libraries for sequencing (Fig. 3).
240
241    *iTru fusion primer experiments*
242    For TruSeq-compatible libraries, we designed and synthetized locus-specific forward fusion
243    primers, which started on the 5' end with the Illumina TruSeq Read1 sequence (5'—
244    ACACTCTTTCCCTACACGACGCTCTTCCGATCT—3') for forward primers or the Illumina
245    TruSeq Read2 sequence (5'—GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT—3') for
246    reverse primers; then included unique five nucleotide (nt) tags (Faircloth & Glenn, 2012) with
247    variable length spacers (0–3 nt) to function as internal indexes (Table 1); and ended with locus-
248    specific primer sequences (Fig. 2; Table 2). To assist with production of fusion primers and
249    reduce errors, we have created and provided Excel spreadsheets (TaggiMatrix; File S1) and a
250    web page (http://baddna.uga.edu/tools-taggi.html). With TaggiMatrix, users can simply input the
251    names and sequences of the locus-specific primers, and all 22 (i.e., 2 non-indexed and 20
252    internally indexed) fusion primers and names are generated automatically. It is important to note
253    that secondary structures or other PCR inhibiting characteristics are not checked by these tools
254    (see Discussion). We then used the locus-specific fusion primers in a primary PCR, followed by
255    a clean-up step and a subsequent PCR with iTru primers from *Adapterama I*. As an example, a
256    general protocol for 16S amplification using TaggiMatrix can be found in File S2.
257        We used this approach for five projects (Table 4), each with slight modifications. First,
258    we used primers targeting *cytochrome-b* to characterize the source of blood meals in kissing
259    bugs; in this project, we first amplified DNA with standard primers, then ligated a y-yoke
260    adapter to these products, and then amplified these products in an iTru PCR (Method 1 in Table
261    3). Second, we used primers targeting several portions of the ITS region, including "flipped"
262    fusion primers, to identify fungal pathogens in tree tissues; in this project, we first amplified
263    DNA with standard primers, then amplified these products with indexed fusion primers, and then
264    amplified these products in an iTru PCR (Method 2 in Table 3). Third, we used primers targeting
265    12S to characterize plethodontid salamander communities from environmental DNA samples; in
266    this project, we first amplified DNA with either internally indexed or non-indexed fusion primers
267    and then amplified these products in an iTru PCR (Methods 4 or 5 in Table 3). Fourth, we used
268    primers targeting two regions of the cyclin-dependent kinase inhibitor *p21* promoter to compare
269    basal DNA methylation of *p21* promoter in two types of human cells; in this project, we first
270    amplified DNA with non-indexed fusion primers and then amplified these products in an iTru
271    PCR (Method 4 in Table 3; Kolli et al., 2019). Fifth, we used primers targeting 16S to
272    characterize bacterial gut microbiomes in wild cotton mice (*Peromyscus leucopus*); in this
273    project, we first amplified DNA with internally indexed fusion primers and then amplified these
274    products in an iTru PCR (Method 5 in Table 3; File S2). Full methods describing the sample
275    collection, DNA extraction, library construction (including detailed descriptions of pooling
276    schemes), and data analysis are detailed in File S3.

277

278 *iNext fusion primer experiments*

279 We generated libraries compatible with Nextera sequencing primers using the same approach as
280 described above for TruSeq-compatible libraries, except that forward fusion primers started with
281 Illumina Nextera Read1 sequence (5'—
282 TCGTCGGCAGCGTCAGATGTGTATAAGAGACAG—3'), and reverse primers started with
283 the Illumina Nextera Read2 sequence (5'—
284 GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAG—3'), and the second PCR used iNext
285 primers from *Adapterama I* (Glenn et al., 2019). We have provided separate sheets within the
286 TaggiMatrix Excel file (File S1) to facilitate the construction of iNext fusion primers.

287 We used this approach in one project. We used primers targeting one chloroplast locus,
288 two mitochondrial loci, and two nuclear loci to perform a fine-scale population genetic analysis
289 of the invasive vine *Wisteria*; in this project, we first amplified DNA with indexed fusion
290 primers and then amplified these products in an iNext PCR (Method 5 in Table 3). Full methods
291 describing the sample collection, DNA extraction, library construction (including detailed
292 descriptions of pooling schemes), and data analysis are included in the File S3.

293

294 *Pooling, Sequencing, Analysis*

295 The methods used for pooling, sequencing and analysis varied among the six projects
296 (File S3), but some general approaches were consistently employed. Amplicon library pools
297 from each of the six projects were pooled with additional samples and sequenced at different
298 times on Illumina MiSeq instruments. The sizes of the amplicons were determined from known
299 sequence targets and verified by agarose gel electrophoresis and known size-standards. We
300 quantified purified amplicon pools using Qubit (Thermo Fisher Scientific Inc, Waltham, MA).
301 We then input the size, concentration, and number of desired reads for amplicon sub-pools and
302 all other samples or sub-pools that would be combined together for a sequencing run into an
303 Excel spreadsheet to calculate the amount of each sub-pool that should be used (an example file
304 of our pooling guide can be found in File S4). We targeted total proportions ranging from <1% to
305 44% of the MiSeq runs (Table 4). We used v.3 600 cycle kits to obtain the longest reads possible
306 for four of the projects and v.2 500 cycle kits for two of the projects, which reduces buy-in costs
307 when shorter reads are sufficient.

308 Following sequencing, results were returned via BaseSpace or from demultiplexing the
309 outer indexes contained in the bcl files using Illumina software (bcl2fastq). Following
310 demultiplexing of the outer indexes, we used Mr. Demuxy
311 (https://github.com/lefeverde/Mr_Demuxy; File S5) or Geneious® to demultiplex samples based
312 on internal indexes.

313 Downstream analyses varied according to the goals of each project and further details are
314 found in File S6. In brief, after demultiplexing, we cleaned raw sequencing data from each
315 project by trimming primers and quality-filtering. Then, we compared sequences from project 1–
316 4 against relevant databases to identify OTUs. For projects 5–6, we mapped reads to appropriate

317    reference sequences. For project 5, we extracted methylation profiles, whereas for project 6, we
318    identified sequencing polymorphisms among genes and individuals. Additional details about
319    each project are presented in Supplemental File S3.
320
321    **Results**
322    We used five methods that take advantage of iTru or iNext indexing primers developed in
323    *Adapterama I* in six exemplar amplicon sequencing projects. These projects illustrate the range
324    of methodological approaches that can be used to overcome challenges of amplicon library
325    preparation and fulfill most of the characteristics of an ideal amplicon library preparation system.
326        In all but one project (Table 4, project 1), we designed fusion primers to generate
327    amplicons that can be amplified by iTru5 and iTru7 (or iNext5 and iNext7) primers to create full-
328    length Illumina TruSeq (or Nextera) libraries. The indexed fusion primers utilize 20 (i.e., 8 + 12)
329    internal identifying sequences with an edit distance ≥ 3 (Table 1) to create up to 96 internally
330    dual-indexed amplicon libraries which were used individually or pooled for additional outer
331    indexing by iTru5 and iTru7 (or iNext5 and iNext7) primers. Sequential PCRs that start with
332    internally indexed primers create quadruple-indexed amplicon libraries that achieve our design
333    goals of cost reduction, facilitation of large-scale multiplexing, increased base-diversity for
334    Illumina sequencing, and maximization of efficiency of library preparation.
335        In our project characterizing the blood meals of kissing bugs (Table 4, project 1), we
336    obtained an average of 116,902 reads for each sample and identified a total of five unique
337    vertebrate species as the source of the blood meals. In our project identifying fungal pathogens in
338    tree tissues (Table 4, project 2), we obtained an average of 436,825 reads per pool (i.e., 96
339    samples) and characterized the diverse fungal communities found in these samples. In our project
340    characterizing plethodontid salamander communities from environmental DNA samples (Table
341    4, project 3), we obtained an average of 163,555 reads for each PCR replicate and identified
342    reads matching 6/7 species expected to be present in the streams. In our project comparing basal
343    DNA methylation of *p21* (Table 4, project 4), we obtained approximately 10,000 reads per
344    sample and detected differences in methylation of CpG sites between embryonic kidney cells and
345    human proximal tubule cell (Kolli et al., 2019). In our project characterizing bacterial gut
346    microbiomes (Table 4, project 5), we rarified to 15,000 quality-filtered reads per sample and
347    identified an average of 3,847 OTUs per sample. In our project focused on the fine-scale
348    population genetic analysis of *Wisteria* (Table 4, project 6), we obtained an average of 1,697
349    reads per sample and discovered little evidence of population structure among samples. Variation
350    in the average number of reads among projects reflects the intentional allocation of reads when
351    pooling with genomic libraries for sequencing; for example, we pooled plates of libraries for the
352    fungal pathogen project in relative quantities intended to generate approximately 4,000 reads per
353    sample. Variation in the number of reads among samples within a given project likely reflects
354    quantification error and variation in input DNA quantity and quality. Full results and associated
355    figures for each project are detailed in File S3.

356     The costs associated with each method vary significantly, and which approach has the
357     lowest cost depends on the number of samples processed (Fig. 4: note axis scales are not linear;
358     Table 5; File S6). Methods 1 and 4 have the lowest buy-in cost, but the cost of library
359     preparations are fixed, rather than decreasing as the number of samples increases. The constant
360     cost per sample is due to the need for individual second round PCRs (e.g., iTru5/7). The other
361     methods allow pooling of samples prior to second round PCR, which reduces costs. Because
362     Method 1, with no use of fusion primers (non-indexed/indexed), has the highest library
363     preparation costs per sample, it quickly becomes the most expensive method, more than doubling
364     the cost of most other methods with as few as 96 samples. Method 4 remains economically
365     reasonable for processing one or two plates of samples but becomes less reasonable as more
366     plates of samples are used. Method 2 is never economically best, but it is sometimes necessary to
367     achieve sufficient amplification to construct the desired libraries. Thus, Method 2 is only viable
368     when the other methods fail. Method 3 has a moderate buy-in cost and the second-lowest cost
369     per sample for large numbers of samples. Also, Method 3 has the lowest cost when ≤11 plates of
370     samples will be processed, though the cost is very similar to Method 5 after ≥2 plates of samples
371     are processed. Method 5 has the second highest buy-in costs, but the lowest costs per sample
372     when large numbers of samples are processed. Method 5 is optimal when >12 plates of samples
373     are processed. Because Methods 3 and 5 are similar in cost after a few plates of samples are
374     processed, other considerations, such as workflow and personnel costs, are likely to drive
375     decisions about the optimal method rather than the costs of reagents.
376
## Discussion
378     In *Adapterama I*, we introduced a general approach to reduce the cost of genomic library
379     preparations for Illumina instruments. Here, we made extensive use of the iNext and iTru
380     primers described in *Adapterama I* and show that these can also be used to facilitate amplicon
381     library construction at reduced cost with increased flexibility. As we did in *Adapterama I*, we
382     focused mostly on iTru to simplify our presentation of the method, but iNext works identically in
383     most situations.
384          Although we focused on Illumina, many of these approaches can be extended to other
385     platforms following the design principles described here (e.g., use primers from sheet
386     ITS_10nt_5'tags in File S1 following Method 3). For platforms that sequence individual
387     molecules (e.g., PacBio and Oxford Nanopore), there is no advantage to variable-length indexes
388     and negligible penalty for longer indexes, but there are significant informatic advantages to
389     equal-length indexes. Thus, for many other platforms, it will be better to use longer indexes of
390     equal length.
391          In general, TaggiMatrix Method 5 achieves our design goals, in that it: 1) uses the
392     universal Illumina sequencing primers; 2) minimizes costs (as little as $2.20 per library, i.e.
393     Method 3 when prepping 1,248 samples in thirteen pools, Figure 4, File S6); 3) minimizes time
394     and equipment needed for library preparations; 4) minimizes buy-in costs through the use of a
395     limited number of fusion primers and universal iTru7 and iTru5 primers; 5) eliminates error-

396   prone ligation steps; 6) allows for > thousands of samples to be pooled and run simultaneously;
397   7) allows users to vary amplicon representation from tiny to large fractions of a sequencing run
398   (up to 91% has been validated for other projects, data not shown); 8) supports creating millions
399   of samples (8 x 12 x 384 x 384 = 14,155,776) that can be tracked and multiplexed through
400   quadruple-indexing. TaggiMatrix Method 3 shares nearly all of these advantages; per sample
401   costs are a few cents more and ligation of a universal stub onto the amplicon pool is maintained.
402          Similar to other *Adapterama* applications, TaggiMatrix offers several methods for
403   combinatorial and hierarchical indexing of samples (Table 3), allowing users to optimize various
404   criteria. For example, different indexes can be used at any combination of the four index
405   positions in the TaggiMatrix library (Fig. 3). By using inner indexes in combination, 20 (8 + 12)
406   indexes can be used to identify 96 (8 x12) samples. By using inner and outer indexes
407   hierarchically, 40 (8 + 12 + 8 + 12) indexes can identify 9216 (8 x 12 x 8 x 12) samples. By
408   using two sets of iTru5 and iTru7 primers, 36,864 (8 x 12 x [8 + 8]x[12 + 12]) samples can be
409   identified. Varying indexes at all index positions is the most economical way to tag samples,
410   especially as the number of samples increases (Table 6). By combining a single set of 20 (8 + 12)
411   fusion primers with the full set of 384 iTru5 and 384 iTru7 primers from *Adapterama I* (Glenn et
412   al., 2019), a total of 14,155,776 (8 x 12 x 384 x 384) samples can be multiplexed.
413          Our methods address the issue of base diversity through the incorporation of indexes with
414   variable-length spacers that allow for diversity at each base position. This strategy is based on
415   independently originating ideas implemented at the Broad Institute, our lab and others, such as
416   the system developed by Fadrosh et al. (2014) where they introduced "heterogeneity spacers" for
417   sequencing amplicons out of phase. Longer spacers (e.g., 0–7 nt) are advantageous over shorter
418   spacers to compensate for longer repeats in the target amplicons. Mononucleotide repeats are
419   particularly problematic in terms of base diversity. Mononucleotide repeats of ≥5 bp will not be
420   addressed by our short spacers (Table 1). Because Illumina reads are of set length, longer spacers
421   decrease the total amount of useful sequence obtained for downstream analyses. Thus, there is a
422   trade-off in how long the heterogeneity spacers should be. Here, we implement a 0–3 nt long
423   heterogeneity spacers, although this could be easily tuned to 0–7 nt for forward primers and 0–11
424   nt for reverse primers, to accommodate any researcher's preferences and mononucleotide repeats
425   known to occur in the target sequences.
426          Our approach does not deal with the limitation of read-length on Illumina platforms. For
427   long amplicons where complete sequencing is desired, it is possible to construct shotgun libraries
428   from the longer amplicons (e.g., using Illumina Nextera XT, Kapa Biosystems Hyper Prep Plus,
429   NEB Ultra II FS or many other commercial kits). The methods used in *Adapterama I* may be
430   helpful in those cases. Such libraries can take advantage of the reduced costs per read on higher
431   capacity instruments. It is also possible to design internal locus-specific fusion primers that
432   recover the entire desired DNA region through independent PCRs. It is important to note,
433   however, that the recent introduction of the PacBio Sequel II along with sequencing chemistry
434   v.6 makes circular consensus sequencing of long amplicons on PacBio an economically
435   reasonable approach. Thus, use of the longer consistent-length indexes noted above to create

436  amplicon pools for PacBio is likely to be increasingly attractive as their platform continues to
437  improve.
438      TaggiMatrix provides an easy way to create indexed fusion primers for convenient
439  ordering at any oligo vendor of your choice. However, the current web page and spreadsheets do
440  not perform quality control of the primer sequences generated. Thus, before ordering, it is
441  important to validate the fusion primers to ensure hairpins, dimers and other secondary structures
442  that inhibit PCR are not created. Several programs exist to validate the primers designed and
443  these should be used before ordering. It is also generally recommended that a small number of
444  fusion primers should be obtained and tested prior to investing large batches of long fusion
445  primers. When deciding on the best method to use (i.e., Methods 1–5), the number of samples,
446  reagent cost, and time available to optimize the primers should be considered (Fig. 5).
447      While developing adapters and primers to make multiple libraries that will be pooled and
448  sequenced, it is important to determine if the primers with different indexes have biased
449  amplification characteristics. This can be accomplished by testing all primers via quantitative
450  PCR using a common template pool to ensure that each primer was synthesized, aliquoted, and
451  reconstituted successfully and has similar amplification efficiency. In practice, however, it will
452  not be economical or necessary to conduct such rigorous quality control for many projects. It is
453  important to note that because sequencing reads are so cheap (~10,000 reads per $1 USD for
454  PE300 reads on a MiSeq), being off by thousands of reads per sample is less expensive than
455  precise quantification, especially when personnel time for such quantification is considered.
456  Thus, it will often be less expensive to subsample reads from overrepresented samples and/or
457  simply redo the small proportion of samples that do not generate a sufficient number of reads.
458  Another common concern with amplicon library preparation methods involving PCR is the
459  introduction of bias due to PCR duplicates. Our method can be modified to incorporate 8N
460  indices similar to how we addressed this issue with RADcap libraries (Hoffberg et al., 2016). It
461  is also possible to use internal N indices of any length desired as molecular identifiers (i.e.,
462  Jabara et al., 2011; Kou et al., 2016). These modifications, in conjunction with long-amplicon
463  sequence on other platforms is worthy of further work.
464

## Conclusions

466  In summary, we demonstrate how several variants of TaggiMatrix solve common challenges for
467  amplicon sequencing on NGS platforms. Our methods can be implemented in projects from a
468  wide array of disciplines such as microbial ecology, molecular systematics, conservation
469  biology, population genetics, and epigenetics, and we encourage others to further develop the
470  tools we provide for solving additional challenges posed by these applications.
471

## Acknowledgements

475  years. We thank John Maerz for his help collecting environmental DNA samples. We thank
476  Bradley Brown for his technical help on Bismark.
477

478  **References**

479  Abarenkov, K., Nilsson, R.H., Larsson, K.H., Alexander, I.J., Eberhardt, U., Erland, S., Høiland,
480      K., Kjøller, R., Larsson, E., Pennanen, T., Sen, R., Taylor, A.F.S., Tedersoo, L., Ursing,
481      B.J., Vrålstad, T., Liimatainen, K., Peintner, U., & Kõljalg, U. (2010). The UNITE
482      database for molecular identification of fungi – recent updates and future perspectives.
483      *New Phytologist*, **186**, 281-285. doi: 10.1111/j.1469-8137.2009.03160.x

484  Ansorge, W.J. (2009). Next-generation DNA sequencing techniques. *Nature Biotechnology,* **25**,
485      195-203. doi: 10.1016/j.nbt.2008.12.009

486  Altschul, S.F., Gish, W., Miller, W., Myers, E.W., & Lipman D.J. (1990). Basic local alignment
487      search tool. *Journal of Molecular Biology*, **215**, 403-410. doi: 10.1016/S0022-
488      2836(05)80360-2

489  Bengtsson-Palme, J., Ryberg, M., Hartmann, M., Branco, S., Wang, Z., Godhe, A., … Nilsson,
490      R.H. (2013). Improved software detection and extraction of ITS1 and ITS2 from
491      ribosomal ITS sequences of fungi and other eukaryotes for analysis of environmental
492      sequencing data. *Methods in Ecology and Evolution*, **4**, 914-919. doi: 10.1111/2041-
493      210X.12073

494  Bentley, G., Higuchi, R., Hoglund, B., Goodrige, D., Sayer, D., Trachtenberg, E.A., & Erlich,
495      H.A. (2009). High-resolution, high-throughput HLA genotyping by next-generation
496      sequencing. *Tissue Antigens*, **74**, 393-403. doi: 10.1111/j.1399-0039.2009.01345.x

497  Binladen, J., Gilbert, M.T.P., Bollback, J.P., Panitz, F., Bendixen, C., … Willerslev, E. (2007).
498      The use of coded PCR primers enables high-throughput sequencing of multiple homolog
499      amplification products by 454 parallel sequencing. *PLoS ONE*, **2**, e197. doi:
500      10.1371/journal.pone.0000197

501  Bybee, S.M., Bracken-Grissom, H., Haynes, B.D., Hermansen, R.A., Byers, R.L., Clement, M.J.,
502      … Crandall, K.A. (2011). Targeted Amplicon Sequencing (TAS): A scalable Next-Gen
503      approach to multilocus, multitaxa phylogenetics. *Genome Biology and Evolution*, **3**,
504      1312-1323. doi: 10.1093/gbe/evr106

505  Caporaso, J.G., Lauber, C.L., Walters, W.A., Berg-Lyons, D., Lozupone, C.A., Turnbaugh, P.J.,
506      … Knight, R. (2011). Global patterns of 16S rRNA diversity at a depth of millions of
507      sequences per sample. *PNAS*, **108**, 4516-4522. doi: 10.1073/pnas.1000080107

508  Craig, D.W., Pearson, J.V., Szelinger, S., Sekar, A., Redman, M., Corneveaux, J.J., …
509      Huentelman, M.J. (2008). Identification of genetic variants using bar-coded multiplex
510      sequencing. *Nature Methods*, **5**, 887-893. doi: 10.1038/nmeth.1251

511  Cronn, R., Knaus, B.J., Liston, A., Maughan, P.J., Parks, M., Syring, J.V., Udall, J. (2012).
512      Targeted enrichment strategies for next-generation plant biology. *American Journal of
513      Botany*, **99**, 291-311. doi: 10.3732/ajb.1100356

514   Cruaud, P., Rasplus, J.Y., Rodriguez, L.J., & Cruaud, A. (2017). High-throughput sequencing of
515         multiple amplicons for barcoding and integrative taxonomy. *Scientific Reports*, **7**, 41948.
516         doi: 10.1038/srep41948
517   DeSantis, T.Z., Hugenholtz, P., Larsen, N., Rojas, M., Brodie, E.L., Keller, K., … Andersen, L.
518         (2006). Greengenes, a chimera-checked 16S rRNA gene database and workbench
519         compatible with ARB. *Applied and Environmental Microbiology,* 72, 5069-5072. doi:
520         10.1128/AEM.03006-05
521   Edgar, R.C. (2010). Search and clustering orders of magnitude faster than BLAST.
522         *Bioinformatics,* **26**, 2460-2461. doi: 10.1093/bioinformatics/btq461
523   Fadrosh, D.W., Ma, B., Gajer, P., Sengamalay, N., Ott, S., Brotman, R.M., & Ravel, J. (2014).
524         An improved dual-indexing approach for multiplexed 16S rRNA gene sequencing on the
525         Illumina MiSeq platform. *Microbiome*, **2**:6. doi: 10.1186/2049-2618-2-6.
526   Faircloth, B.C., & Glenn, T.C. (2012). Not all sequence tags are created equal: Designing and
527         validating sequence identification tags robust to indels. *PLoS ONE*, **7**, e42543. doi:
528         10.1371/journal.pone.0042543
529   Glenn, T.C. (2011). Field guide to next-generation DNA sequencers. *Molecular Ecology*
530         *Resources,* **11**, 759-769. doi: 10.1111/j.1755-0998.2011.03024.x
531   Glenn, T.C., Nilsen, R.A., Kieran, T.J., Sanders, J.G., Bayona-Vásquez, N.J., Finger, J.W. Jr., …
532         Faircloth, B.C. (2019). Adapterama I: Universal stubs and primers for 384 unique dual-
533         indexed or 147,456 combinatorially-indexed Illumina libraries (iTru & iNext). *BioRxiv,*
534         doi: 10.1101/049114.
535   Goodwin, S., McPherson, J.D., McCombie, W.R. (2016). Coming of age: ten years of next-
536         generation sequencing technologies. *Nature Reviews, Genetics*, **17**, 333-351. doi:
537         10.1038/nrg.2016.49
538   Hoffberg, S.L., Kieran, T.J., Catchen, J.M., Devault, A., Faircloth, B.C., Mauricio, R., Glenn,
539         T.C. (2016). RADcap: Sequence capture of dual-digest RADseq libraries with
540         identifiable duplicates and reduced missing data. *Molecular Ecology Resources* **16**, 1264-
541         1278. doi: 10.1111/1755-0998.12566
542   Illumina. (2018a). Nextera XT DNA library prep kit: Reference guide. Illumina Proprietary
543         Document # 15031942v03, February 2018.
544         https://support.illumina.com/content/dam/illumina-
545         support/documents/documentation/chemistry_documentation/samplepreps_nextera/nexter
546         a-xt/nextera-xt-library-prep-reference-guide-15031942-03.pdf, accessed 2 April, 2019.
547   Illumina. (2018b). Indexed sequencing overview guide. Illumina Proprietary Document
548         #15057455v04, February 2018. http://support.illumina.com/content/dam/illumina-
549         support/documents/documentation/system_documentation/miseq/indexed-sequencing-
550         overview-guide-15057455-04.pdf, accessed 5 October 2018.
551   Jabara, C.B., Jones, C.D., Roach, J., Anderson, J.A., Swanstrom, R. (2011). Accurate sampling
552         and deep sequencing of the HIV-1 protease gene using a Primer ID. *Proc Natl Acad Sci*
553         *USA*, **108(50)**, 20166– 20171. doi: 10.1073/pnas.1110064108

554  Kieran, T.J., Gottdenker, N.L., Varian, C.P., Saldaña, A., Means, N., Owens, D., … Glenn, T.
555       (2017). Bloodmeal source characterization using Illumina sequencing in the Chagas
556       Disease vector *Rhodnius pallescens* (Hemiptera: Reduviidae) in Panama. *Journal of*
557       *Medical Entomology,* **54(6)**, 1786-1789. doi: 10.1093/jme/tjx170

558  Klindworth, A., Pruesse, E., Schweer, T., Peplies, J., Quast, C., Horn, M., & Glöckne, O. (2013).
559       Evaluation of general 16S ribosomal RNA gene PCR primers for classical and next-
560       generation sequencing-based diversity studies. *Nucleic Acids Research,* **41**, e1. doi:
561       10.1093/nar/gks808

562  Kolli, R.T., Glenn, T.C., Brown, B.T., Kaur, S.P., Barnett, L.M., Lash, L.H., Cummings, B.
563       (2019). Bromate-induced changes in *p21* DNA methylation and histone acetylation in
564       renal cells. *Toxicological Sciences*, **168**, 460-473. doi: 10.1093/toxsci/kfz016

565  Kou, R., Lam, H., Duan, H., Ye, L., Jongkam, N., Chen, W., … Shihong, Li. (2016). Benefits
566       and challenges with applying unique molecular identifiers in Next Generation
567       Sequencing to detect low frequency mutations. *PLoS ONE*, **11(1**): e0146638. doi:
568       10.1371/journal.pone.0146638

569  Langmead, B., & Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie 2. *Nature*
570       *Methods*, **9**, 357-359. doi: 10.1038/nmeth.1923

571  Magoč, T., & Salzberg, S. (2011). FLASH: fast length adjustment of short reads to improve
572       genome assemblies. *Bioinformatics*, **27,** 2957-2963. doi: 10.1093/bioinformatics/btr507

573  Meirmans, P.G., & Van Tienderen, P.H. (2004). GENOTYPE and GENODIVE: two programs
574       for the analysis of genetic diversity of asexual organisms. *Molecular Ecology Notes,* **4**,
575       792-794. doi: 10.1111/j.1471-8286.2004.00770.x

576  Mitra, A., Skrzypczak, M., Ginalski, K., & Rowicka, M. (2015). Strategies for achieving high
577       sequencing accuracy for low diversity samples and avoiding sample bleeding using
578       Illumina platform. *PLoS ONE*. **10**, e0120520. doi: 10.1371/journal.pone.0120520

579  Noireau, F., Abad-Franch, F., Valente, S.A., Dias-Lima, A., Lopes, C.M., Cunha, V., … Jurberg,
580       J. (2002). Trapping Triatominae in silvatic habitats. *Mem Inst Oswaldo Cruz*, **97**, 61-63.
581       doi: 10.1590/S0074-02762002000100009

582  Parson, W., Pegoraro, K., Niederstätter, H., Föger, M., & Steinlechner, M. (2000). Species
583       identification by means of the cytochrome b gene. International Journal of Legal
584       Medicine, **114**, 23-28. doi: 10.1007/s004140000

585  Pierson, T.W., McKee, A.M., Spear, S.F., Maerz, J.C., Camp, C.D., & Glenn, T.C. (2016).
586       Detection of an enigmatic plethodontid salamander using environmental DNA. *Copeia*,
587       **2016(1)**, 78-82. doi: 10.1643/CH-14-202.

588  Quail, M.A., Smith, M., Coupland, P., Otto, T.D., Harris, S.R., Connor T.R., Bertoni, A.,
589       Swerdlow, H.P., & Gu, Y. (2012). A tale of three next generation sequencing platforms:
590       comparison of Ion Torrent, Pacific Biosciences and Illumina MiSeq sequencers. *BMC*
591       *Genomics*, **13**: 341-353. doi: 10.1186/1471-2164-13-341

592  Shokralla, S., Gibson, J.F., Nikbakht, H., Janzen, D.H., Hallwachs, W., & Hajibabae, M. (2014).
593       Next-generation DNA barcoding: using next-generation sequencing to enhance and

594        accelerate DNA barcode capture from single specimens. *Molecular Ecology Resources*,

595        **14**, 892-901. doi: 10.1111/1755-0998.12236

596  Toju, H., Tanabe, A.S., Yamamoto, S., & Sato, H. (2012). High-Coverage ITS Primers for the

597        DNA-Based Identification of Ascomycetes and Basidiomycetes in Environmental

598        Samples. *PLoS ONE*. **7**, e40863. doi: 10.1371/journal.pone.0040863

599  Trusty, J.L., Goertzen, L.R., Zipperer, W.C., & Lockaby, B.G. (2007a). Invasive Wisteria in the

600        Southeastern United States: genetic diversity, hybridization and the role of urban centers.

601        *Urban Ecosystems*, **10**, 379-395. doi: 10.1007/s11252-007-0030-y

602  Trusty, J.L., Lockaby, B.G., Zipperer, W.C., & Goertzen, L.R. (2007b). Identity of naturalised

603        exotic Wisteria (Fabaceae) in the south-eastern United States. *Weed Research*, **47**, 479-

604        487. doi: 10.1111/j.1365-3180.2007.00587.x

605  Trusty, J.L., Lockaby, B.G., Zipperer, W.C., & Goertzen, L.R. (2008). Horticulture, hybrid

606        cultivars and exotic plant invasion: a case study of Wisteria (Fabaceae). *Botanical*

607        *Journal of the Linnean Society*, **158**, 593-601. doi: 10.1111/j.1095-8339.2008.00908.x

608  White, T.J., Bruns, T., Lee, S., & Taylor, J. (1990). Amplification and direct sequencing of

609        fungal ribosomal RNA genes for phylogenetics. In MA Innis, DH Gelfand, JJ Sninky &

610        TJ White (Eds.), *PCR Protocols: A guide to methods and applications* (pp. 315-322) San

611        Diego, USA: Academic Press.

**Table 1**
**Internal identifying index sequences.**
All indexes have an edit distance of ≥ 3.  Upper case letters are the indexes; lower case letters add length variation to facilitate sequence diversity at each base position of amplicon pools (see text for details).  For Illumina MiSeq and HiSeq models ≤2500, adenosine and cytosine are in the red detection channel, whereas guanine and thymine are in the green channel.  Indexes and spacers have balanced red and green representation at each base position within each group of four indexes (i.e., count 1–4, 5–8, 9–12, 13–16, and 17–20).

| Index count | Index Label | Sequence | Length |
|---|---|---|---|
| 1 | A | GGTAC | 5 |
| 2 | B | cAACAC | 6 |
| 3 | C | atCGGTT | 7 |
| 4 | D | tcgGTCAA | 8 |
| 5 | E | AAGCG | 5 |
| 6 | F | gCCACA | 6 |
| 7 | G | ctGGATG | 7 |
| 8 | H | tgaTTGAC | 8 |
| 9 | 1 | AGGAA | 5 |
| 10 | 2 | gAGTGG | 6 |
| 11 | 3 | ccACGTC | 7 |
| 12 | 4 | ttcTCAGC | 8 |
| 13 | 5 | CTAGG | 5 |
| 14 | 6 | tGCTTA | 6 |
| 15 | 7 | gcGAAGT | 7 |
| 16 | 8 | aatCCTAT | 8 |
| 17 | 9 | ATCTG | 5 |
| 18 | 10 | gAGACT | 6 |
| 19 | 11 | cgATTCC | 7 |
| 20 | 12 | tctCAATC | 8 |

**Table 2**

**Primer pairs used in the example projects presented.**

Project, target locus, forward and reverse primer names and sequences, as well as the sources of the primer sequences are shown.

| Project | Target Locus | Forward Primer | Reverse Primer |
|---|---|---|---|
| Kissing Bug[1] | cyt-b | L14816: CCATCCAACATCTCAGCATGATGAAA | H15173: CCCCTCAGAATGATATTTGTCCTCA |
| Pathogenic Fungi[2,3] | ITS | ITS1-F_KYO2: TAGAGGAAGTAAAAGTCGTAA | ITS2_KY02: TTYRCTRCGTTCTTCATC |
| | | ITS3-KYO2: AHCGATGAAGAACRYAG | ITS4: TCCTCCGCTTATTGATATGC |
| | | ITS1-F_KYO2: TAGAGGAAGTAAAAGTCGTAA | ITS4: TCCTCCGCTTATTGATATGC |
| Salamander eDNA | 12S | Pleth_12S_F: AAAAAAGTCAGGTCAAGG | Pleth_12S_R: GGTGACGGGCGGTGTGTG |
| Bacterial Community[4,5] | 16S | Bact-0341-b-S-17: CCTACGGGNGGCWGCAG | S-D-Bact-0785-a-A-21: GACTACHVGGGTATCTAATCC |
| | 16S | 515F: GTGCCAGCMGCCGCGGTAA | 806R: GGACTACHVGGGTWTCTAAT |
| Methylation[6] | *p21-TSS* | hp21-TSS F: ATAGTGTTGTGTTTTTTTGGAGAGTG | hp21-TSS R: ACAACTACTCACACCTCAACTAAC |
| | *SIE-1* | hp21-SIE1 F: TTTTTTGAGTTTTAGTTTTTTTAGTAGTGT | hp21-SIE1 R: AACCAAAATAATTTTTCAATCCC |
| *Wisteria*[7,8,9] | nr824 | w898-824F: CATGTTGCATTCAATCTTGG | w898-824R: GCCTCCATACAAGTTAGTTG |
| | nr997 | w843-997F: GAATCAACGCTGAACGTT | w843-997AluR: GGTTCAATTTATTGATGTG |
| | trnL; trnL/F | WistmLF: AGTTGACGACATTTCCTTAC | WistmLR: GGAGTGAATGGTTTGATCAATG |
| | nad4 | NAD4RSF1: CTACTAGACTACTAGAGGT | NAD4RSRl: GTTTGGCAACAAGCAAACG |
| | cyt-b | COBRSF1: CATATTGACTTTCTCTCGCC | COBRSR1: GAATAGGATGACTCAGCGTC |

[1] Parson et al. 2000; [2] Toju *et al.* 2012; [3] White et al. 1990; [4] Klindworth et al. 2013; [5] Caporaso et al. 2012; [6] Koli et al. 2018; [7] Trusty et al. 2007a; [8] Trusty et al. 2007b; [9] Trusty et al. 2008

**Table 3**
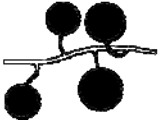**General strategies for producing and indexing amplicon libraries for Illumina sequencing.**
These examples use iTru primers, but as mentioned in the text, this can be implemented instead with iNext primers. Method 5 is illustrated below, but we are not including any dataset in the present manuscript that has implemented it (see Discussion). Note: this table does not include "flipped" primers.

| Method 1 | Method 2 | Method 3 | Method 4 | Method 5 |
|---|---|---|---|---|
| Standard primers | Standard primers | Indexed primers | Fusion primers | Indexed fusion primers |
| ↓ | ↓ | ↓ | ↓ | ↓ |
| PCR | PCR | PCR | PCR | PCR |
|  | ↓ |  |  |  |
| ↓ | Indexed fusion primers | [ Pool ] |  | [ Pool ] |
|  | ↓ | ↓ | ↓ |  |
| Y-yoke | PCR | Y-yoke |  | ↓ |
| ↓ | [Pool] | ↓ |  |  |
| iTru PCR | ↓ | iTru PCR | iTru PCR | iTru PCR |
|  | iTru PCR |  |  |  |
| ↓ | ↓ | ↓ | ↓ | ↓ |
| Completed library | Completed library | Completed library | Completed library | Completed library |

| Method 1 | Method 2 | Method 3 | Method 4 | Method 5 | |
|---|---|---|---|---|---|
| - | + | + | - | + | Base diversity in reads |
| - | + | + | - | + | Poolable to reduce library preparation costs |
| 2 | 20 | 20 | 2 | 20 | Number of primers |
| 192 | 193 | 97 | 192 | 97 | Minimum number of PCRs for 96 samples |
| - | - | + | - | + | PCR bias varies among |
| Low | Low | Med | Med | High | Optimization difficulty |
| Low | High | Med | Med | High | Relative primer cost |
| High | Med | Med | Med | Low | Relative library preparation cost |

**Table 4**

**Detailed information for example projects presented to validate our approach.**

Summarized information for all example projects used to demonstrate Taggimatrix. The "Method" column refers to methods in Table 3; the "Target Reads" column cites the approximate number of reads per pool (i.e., not per individual sample) we targeted when pooling samples with other libraries. Note that these data were generated on many independent MiSeq runs. The kissing bug image is from Joseph Hughes (https://creativecommons.org/licenses/by-nc-sa/3.0/), and all other images are from PhyloPic 2.0 (Public Domain Dedication 1.0).

| # | Organisms | Project Goal | Target Loci | Library Type | Method | Pool Name | Target Reads | Actual Reads | Summary |
|---|---|---|---|---|---|---|---|---|---|
| 1 |  | Diet analysis | cyt-b | iTru | 1 | N/A | 100k (< 1%) | 916k | Identified five vertebrate sources of blood meals. |
| 2 |  | Fungal identification | Full-ITS1 (standard & "flipped") | iTru | 2 | Homokaryon | 400k (2.7%) | 515k | Identified the primary fungal OTU from each culture |
| | | | | | | Het.multispore | 400k (2.7%) | 619k | |
| | | | | | | Het. Tissue | 400k (2.7%) | 444k | |
| | | | Full-ITS2 (standard & "flipped") | iTru | 2 | Homokaryon | 400k (2.7%) | 268k | |
| | | | | | | Het.multispore | 400k (2.7%) | 310k | |
| | | | | | | Het. Tissue | 400k (2.7%) | 257k | |
| | | | Incomplete-ITS1&ITS2 | iTru | 2 | Homokaryon | 400k (2.7%) | 460k | |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | (standard & "flipped") | | | Het.multispore | 400k (2.7%) | 579k | |
| | | | | | | Het. Tissue | 400k (2.7%) | 514k | |
| 3 | | Environmental DNA | 12S | iTru | 4 & 5 | Reference samples | 10k (< 1%) | 8k | Detected 6/7 species of salamander expected in community |
| | | | | | | eDNA samples | 12M (48%) | 4.4M | |
| 4 | | Methylation | *p21-TSS SIE-1* | iTru | 4 | N/A | 40k (0.3%) | 121k | Compared methylation patterns between cell types |
| 5 | | Microbiome | 16S | iTru | 5 | Ash Basin | 1.5M (6%) | 3.8M | Detected 90,862 bacterial OTUs |
| | | | | | | Pond B | 1.5M (6%) | 2.8M | |
| | | | | | | Tim's Branch | 1.5M (6%) | 0.7M | |
| | | | | | | Upper Three Runs | 1.5M (6%) | 2.9M | |
| 6 | | Population genetics | nr824 nr997 trnL; trnL/F nad4 cyt-b | iNext | 5 | N/A | 150k (1.3%) | 79k | Demonstrated mixed ancestry and no population structure in an introduced population |

**Table 5**

**Oligos and iTru buy-in, and library prep costs among methods.**

Costs associated to the implementation of the different methods. In segment **a)** we present buy-in cost of oligos and iTru primers and cost per sample of library prep which consists of both, fixed and variable costs depending on pooling at early stages. Segment **b)** is the cost of library prep (no considering primers/adapters) per sample given a number of samples. Segment **c)** is the total experimental cost of primers/adapters and library prep according to the number of samples in the experiment, the first section is in term of number of samples, the second section is in terms of plates, each plate consisting of 96 samples. Cost for iTru are calculated list prices of aliquots from baddna.uga.edu. Costs for 'oligos' are calculated using list prices from Integrated DNA Technologies (IDT; Coralville, IA). Other costs are from listed prices from various vendors by Jan 2019. Please view File S1 and S6 for additional details on price calculations and also to review total prices of experiment given a number of samples.

**a)**

|  | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| iTru buy-in | $500 | $500 | $500 | $500 | $500 |
| Oligo buy-in | $103 | $460 | $290 | $40 | $445 |
| Library Cost per sample | $18.86 | variable | variable | $4.44 | variable |
| Fixed cost | $18.86 | $3.12 | $1.39 | $4.44 | $1.39 |
| Variable cost | - | $4.07 | $17.52 | - | $4.07 |

**b)**

| | # samples | Library Cost per Sample for the given # of samples | | | | |
|---|---|---|---|---|---|---|
| | 1 | $18.86 | $7.19 | $18.91 | $4.44 | $5.46 |
| | 2 | $18.86 | $5.16 | $10.15 | $4.44 | $3.43 |
| | 8 | $18.86 | $3.63 | $3.58 | $4.44 | $1.90 |
| | 12 | $18.86 | $3.46 | $2.85 | $4.44 | $1.73 |
| | 24 | $18.86 | $3.29 | $2.12 | $4.44 | $1.56 |
| | 48 | $18.86 | $3.20 | $1.75 | $4.44 | $1.47 |
| | 96 | $18.86 | $3.16 | $1.57 | $4.44 | $1.43 |

**c)**

| | | Total Experiment Cost for given # of samples or plates (96 samples per plate) | | | | |
|---|---|---|---|---|---|---|
| # samples | 1 | $621.86 | $967.19 | $808.91 | $544.44 | $950.46 |
| | 2 | $640.72 | $970.31 | $810.30 | $548.87 | $951.85 |
| | 8 | $753.87 | $989.03 | $818.64 | $575.48 | $960.19 |
| | 12 | $829.31 | $1,001.50 | $824.20 | $593.22 | $965.75 |
| | 24 | $1,055.62 | $1,038.94 | $840.87 | $646.45 | $982.43 |
| | 48 | $1,508.24 | $1,113.80 | $874.23 | $752.90 | $1,015.78 |
| | 96 | $2,413.48 | $1,263.53 | $940.94 | $965.80 | $1,082.49 |
| # plates | 2 | $4,223.96 | $1,567.06 | $1,091.87 | $1,391.60 | $1,219.98 |
| | 3 | $6,034.44 | $1,870.59 | $1,242.81 | $1,817.40 | $1,357.47 |
| | 4 | $7,844.92 | $2,174.12 | $1,393.74 | $2,243.20 | $1,494.95 |
| | 5 | $9,655.40 | $2,477.66 | $1,544.68 | $2,669.00 | $1,632.44 |

**Note: These will be added individually to PeerJ with each file upload.** *Don't* **include "Figure 1"; just add the title and description separately. Titles are in bold and descriptions are in plain font.**

Figure 1
**High throughput workflow to create and multiplex TaggiMatrix libraries**
The components of the quadrupled-indexed amplicon Libraries. A specific DNA region is amplified using fusion and tagged locus-specific primers, also known as "indexed fusion primers", to produce a fusion amplicon. Then iTru adapters are ligated using Y-yolk adapters or incorporated using limited cycle PCR with i5 and i7 indexed primers to make the complete double stranded DNA library. Internal indexes and outer i5/i7 indexes are represented as well as the set of primers used.

Figure 2
**Examples of possible primer types (Table 3), including "flipped" fusion primers**
Elements in the box are combined to form each of these various primer types, shown below the box. Standard locus-specific primer sequences are indicated by the letter "N", in uppercase the forward primer and lowercase the reverse primer. Green and red nucleotide bases refer to unique index sequences. Blue and pink sequences are Read1 and Read 2 fusion sequences, respectively.

Figure 3
**Sequencing reads that can be obtained from dual-indexed paired-end reads.**
a) Illustration of a double-stranded DNA molecule from a full-length amplicon library (i.e., following the limited-cycle round of PCR). Horizontal arrowheads indicate the 3' ends. Labels on the double-stranded DNA indicate the function of each section, with shading to help indicate boundaries. b) Scheme of the four separate primers used for the four sequencing reactions that occur in paired-end dual-indexed sequencing and the reads that each primer produces (number in the circle). The four sequencing primers are added one at a time in the following order – Read1, Index Read1, Index Read2, and Read2. Vertical height indicates this order (top primer added first). 3A and 3B correspond to workflow A (NovaSeq™ 6000, MiSeq™, HiSeq 2500, and HiSeq 2000) and workflow B (iSeq™ 100, MiniSeq™, NextSeq™, HiSeq X, HiSeq 4000, and HiSeq 3000), respectively, of dual-indexed workflows on paired-end flow cells (Illumina 2018).

Figure 4
**Total cost of experiments across the five methods given a number of samples.**
Line plot of price of each method according to the number of samples. The starting point in the X-axis (x=0) represents the buy-in cost of oligos.

Figure 5

**Decision tree to select the best fitting method according to the experiment goals and budget.**

Guide of choices to drive an informed decision over the method for amplicon sequencing that may be fit the best for your lab/research/experiment goals.

Supplementary Figure S1

**Diagram of full-length amplicon TaggiMatrix library product**

Double stranded amplicon library product after implementation of TaggiMatrix. Indication tags and indexes incorporated through the use of Fusion primers and iTru/iNext primers, respectively.

Supplementary Figure S2

**Detailed illustration of the components on one of the possible designs (Method 5) to construct TaggiMatrix amplicon libraries**

First, locus specific fusion primers with tags are used to amplify the target DNA region. From this step pooling is possible thanks to the presence of indexes. Then library amplification with the use of iTru universal primers with indexes that allows pool labeling and incorporation of Illumina platform oligos (P5 and P7).

Supplementary File S1

**TaggiMatrix spreadsheet**

Excel spreadsheet demonstrating the step-by-step process to create indexed fusion primers with TaggiMatrix. The first sheet (Introduction) is an introductory explanation of how the document works. The second, third, and fourth sheets (…iTru_Fusions) are examples of the creation of indexed fusion primers for 16S, cyt-b and COI universal primers, respectively. The fifth sheet (iNext_&_iTru_Primers) is a list of the universal primer sequences and prices. The sixth and seventh sheets (…Order_Sheet) are examples of how to fill the order form to fill plates with primer sets. The eighth sheet (PCR_Setup) indicates how to combinatorically layout the primers for a 96-well plate. The ninth, tenth, and eleventh sheets (…Tags…) list the index sequences that are incorporated to the fusion primers, their spacers, and examples.

Supplementary File S2

**TaggiMatrix protocol for 16S amplicon library prep**

Step-by-step library construction for 16S libraries with indexed fusion primers**.**

Supplementary File S3

**Supplementary methods and results for TaggiMatrix example datasets**

A detailed guide through the methods, results, and discussion of sequence analyses from TaggiMatrix data generated for each example dataset presented in this manuscript.

Supplementary File S4

**TaggiMatrix video: what is happening inside the tube?**

This presentation demonstrates the key features of TaggiMatrix, including how the combinatorial indexing is performed in a plate.


Supplementary File S5

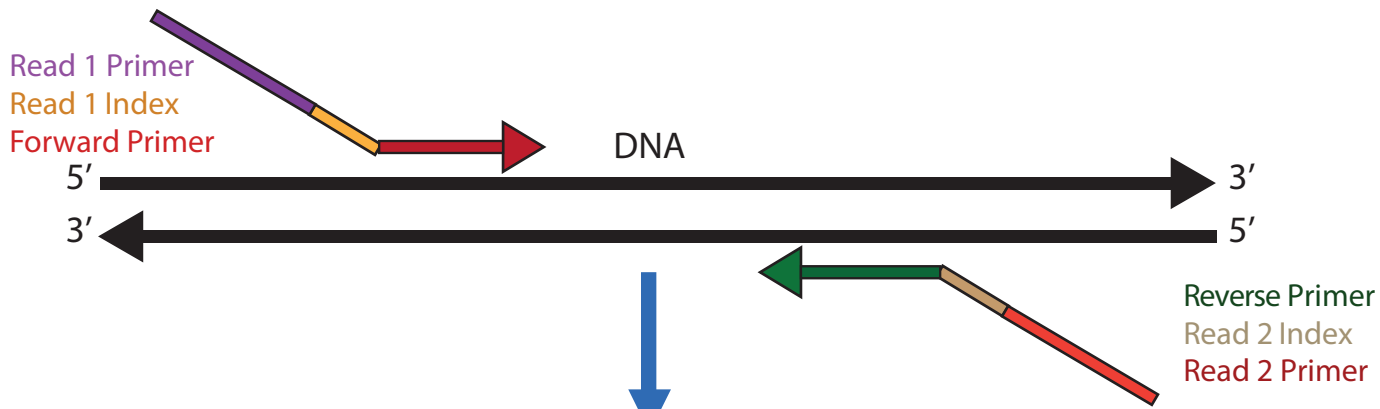**Demultiplexing Internal Indexes Using Mr. Demuxy**

Guide of how to run Mr. Demuxy to demultiplex using internal indexes amplicon data in fastq format.


Supplementary File S6

**Price calculator among methods presented for amplicon sequencing**
Excel spreadsheet with calculations of oligos and reagents costs for library prep among the five methods presented in *Adapterama II*. Users can modify values according to their particular vendors, number of samples, and number of pools, to have an estimate of the price per sample and the price of the experiment.

**Template DNA + Locus-specific Fusion Primers**

Read 1 Primer
Read 1 Index
Forward Primer

DNA

5′ ———————————————————→ 3′
3′ ←——————————————————— 5′

Reverse Primer
Read 2 Index
Read 2 Primer

**Fusion Amplicon**

DNA

**Limited cycle PCR**

Read 1 internal index

Read 2 internal index

i5 primer

i7 primer

i5 index

i7 index

**Double stranded DNA library**

p5

DNA

p7

## Locus-specific primers (Standard Primers)

Forward      NNNNNNNNNNNNNNNNNNN

Reverse      nnnnnnnnnnnnnnnnnnn

### Index Sequence

GGTAC

AGGAA

## Universal 5' TruSeqHT

iTru_R1_5'    ACACTCTTTCCCTACACGACGCTCTTCCGATCT

iTru_R2_5'    GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCT

## Standard Primers with Internal Indexes (Indexed Primers)

Forward      GGTACNNNNNNNNNNNNNNNNNNN

Reverse      GGAAnnnnnnnnnnnnnnnnnnn

## iTru Fusion Primers without Internal Indexes (Fusion Primers)
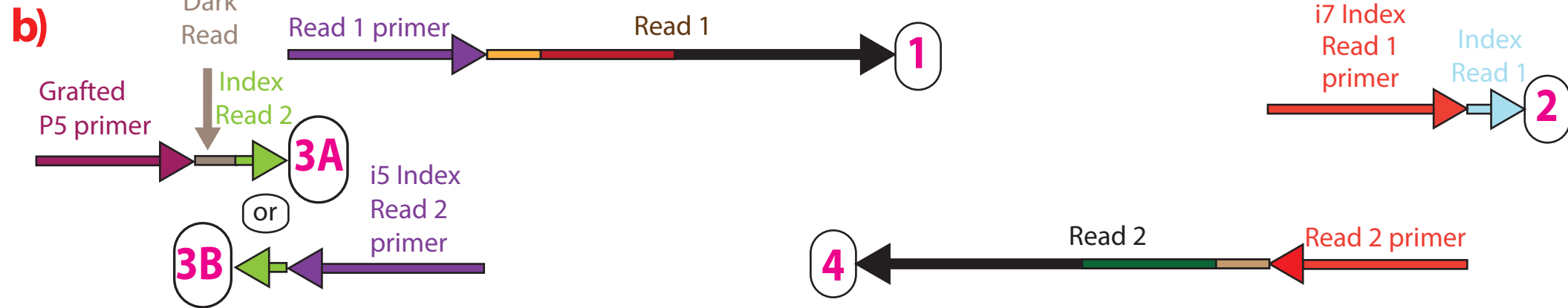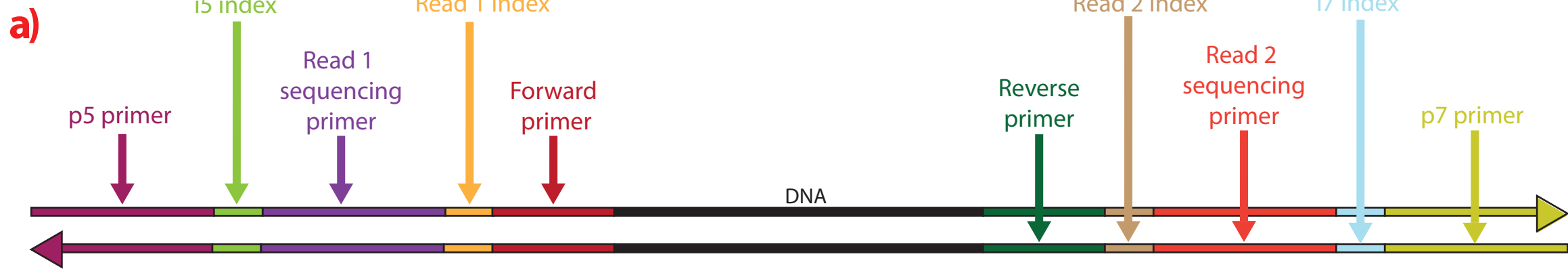
Forward      ACACTCTTTCCCTACACGACGCTCTTCCGATCTNNNNNNNNNNNNNNNNNNN

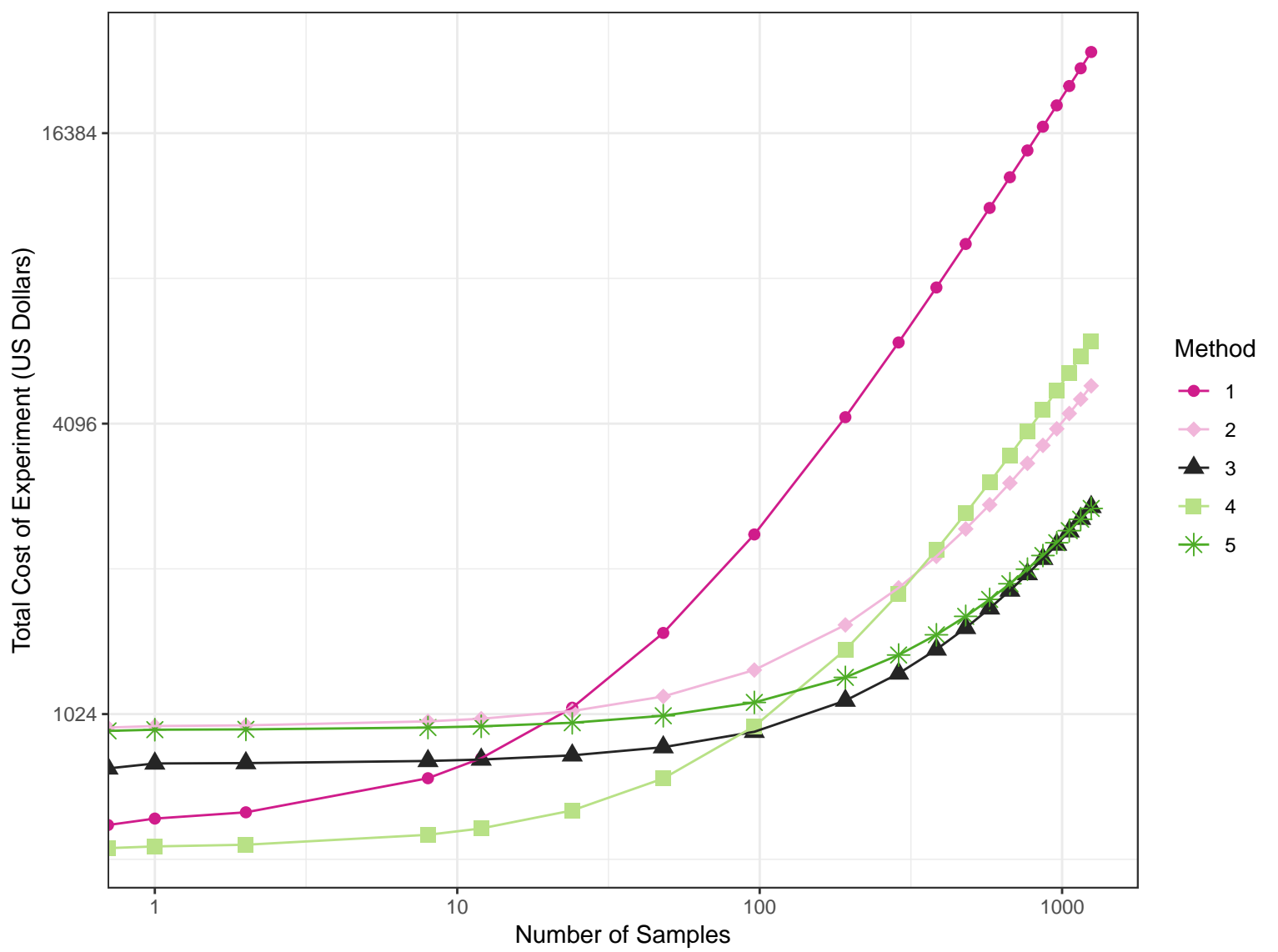Reverse      GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCTnnnnnnnnnnnnnnnnnnn

## iTru Fusion Primers with Internal Indexes (Indexed Fusion Primers)
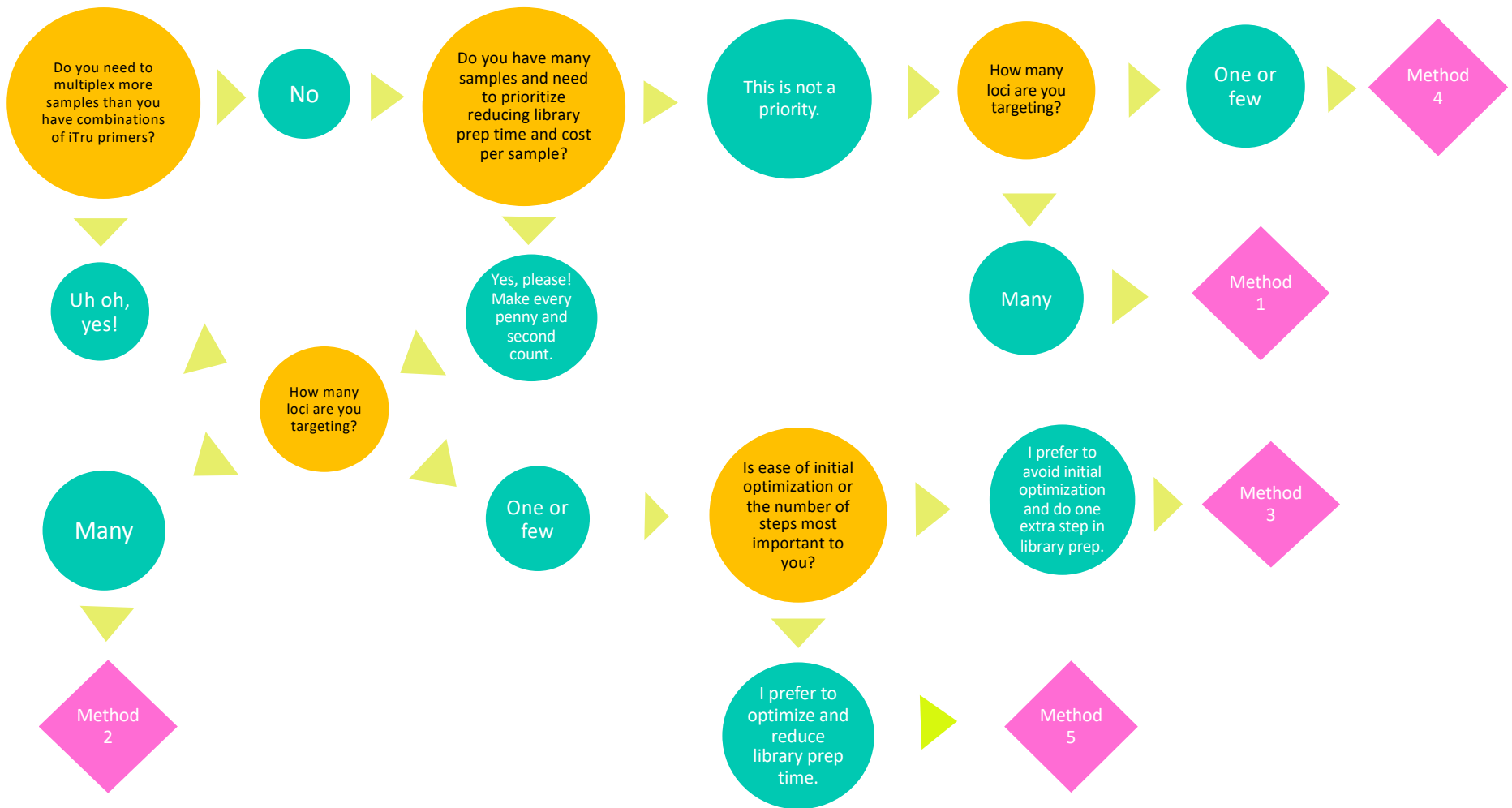
Forward      ACACTCTTTCCCTACACGACGCTCTTCCGATCTGGTACNNNNNNNNNNNNNNNNNNN

Reverse      GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCTAGGAAnnnnnnnnnnnnnnnnnnn

## "Flipped" iTru Fusion Primers with Internal Indexes

Reverse      ACACTCTTTCCCTACACGACGCTCTTCCGATCTGGTACnnnnnnnnnnnnnnnnnnn

Forward      GTGACTGGAGTTCAGACGTGTGCTCTTCCGATCTAGGAANNNNNNNNNNNNNNNNNNN

Do you need to multiplex more samples than you have combinations of iTru primers?

No

Do you have many samples and need to prioritize reducing library prep time and cost per sample?

This is not a priority.

How many loci are you targeting?

One or few

**Method 4**

Uh oh, yes!

How many loci are you targeting?

Yes, please! Make every penny and second count.

Many

**Method 1**

Many

One or few

Is ease of initial optimization or the number of steps most important to you?

I prefer to avoid initial optimization and do one extra step in library prep.

**Method 3**

**Method 2**

I prefer to optimize and reduce library prep time.

**Method 5**

The authors declare competing interests. The EHS DNA lab provide oligonucleotide aliquots and library preparation services at cost, including some oligonucleotides and services that make use of the adapters and primers presented in this manuscript (baddna.uga.edu). The information we present allows all researchers to synthesize the oligonucleotides at any vendor of their choice, follow or modify the library preparation techniques we have included, and freely publish results simply with proper attribution of this paper and Illumina®™. Any use of trade, firm, or product names is for descriptive purposes only and does not imply endorsement by the U.S. Government.