

1 *In silico* assessment of primers for eDNA studies using PrimerTree and application to characterize the
2 biodiversity surrounding the Cuyahoga River

3

4

5 Cannon MV¹, Hester J^{1,2}, Shalkhauser A¹, Chan ER^{1,3}, Logue K¹, Small ST^{1,4} and Serre D^{1*}

6

7 ¹Genomic Medicine Institute, Cleveland Clinic, 9500 Euclid Ave., Cleveland OH 44195, USA

8 ²current address: RStudio, 250 Northern Ave, Boston, MA 02210, USA

9 ³current address: Institute for Computational Biology, Case Western Reserve University, 10900 Euclid
10 Ave., Cleveland OH 44106, USA

11 ⁴Department of Global Health and Diseases, Case Western Reserve University, 10900 Euclid Ave.,
12 Cleveland OH 44106, USA

13

14 * To whom correspondence should be addressed. Tel: (216) 444-0676; Fax: 1 216 636-0009; Email:

15 serred@ccf.org

16 David Serre
17 Genomic Medicine Institute
18 Cleveland Clinic
19 9500 Euclid Ave.
20 Cleveland OH 44195, USA

21

22

23

24

25

1 Analysis of environmental DNA (eDNA) enables the detection of species of interest from water and soil
2 samples, typically using species-specific PCR. Here, we describe a method to characterize the
3 biodiversity of a given environment by amplifying eDNA using primer pairs targeting a wide range of
4 taxa and high-throughput sequencing for species identification. We tested this approach on 91 water
5 samples of 40 mL collected along the Cuyahoga River (Ohio, USA). We amplified eDNA using 12 primer
6 pairs targeting mammals, fish, amphibians, birds, bryophytes, arthropods, copepods, plants and several
7 microorganism taxa and sequenced all PCR products simultaneously by high-throughput sequencing.
8 Overall, we identified DNA sequences from 15 species of fish, 17 species of mammals, 8 species of birds,
9 15 species of arthropods, one turtle and one salamander. Interestingly, in addition to aquatic and semi-
10 aquatic animals, we identified DNA from terrestrial species that live near the Cuyahoga River. We also
11 identified DNA from one Asian carp species invasive to the Great Lakes but that had not been previously
12 reported in the Cuyahoga River. Our study shows that analysis of eDNA extracted from small water
13 samples using wide-range PCR amplification combined with high-throughput sequencing can provide a
14 broad perspective on biological diversity.

15

16 **Introduction**

17 Environmental samples, such as river and pond water or soil, contain a complex mixture of DNA
18 molecules originating from intact microorganisms, feces, mucous, gametes, shed tissues or decaying
19 parts from organisms living in or near the sampling site¹. DNA extracted from these samples (often
20 referred to as environmental DNA or eDNA) can be characterized globally by shotgun sequencing all
21 DNA present within a sample. This approach (referred to as metagenomics) provides a wealth of
22 information, not only about the identity of the species present in the environment, but also about the
23 gene content of the sequenced organisms which can highlight interesting biological features^{2,3}.

1 However, this approach requires extensive sequencing to capture the biological complexity of each
2 sample and is therefore expensive to analyze many samples. Alternatively, DNA extracted from
3 environmental samples can be analyzed by PCR targeting a carefully selected locus. This approach has
4 been widely applied to ecological studies⁴⁻⁶, conservation^{7,8} or to identify the presence of invasive
5 species⁹. This approach is particularly cost-efficient since all reads generated are informative with regard
6 to species identification, enabling detection of even rare organisms. Such eDNA studies typically focus
7 on the analysis of a single macroorganism species and therefore rely on species-specific amplification of
8 DNA (i.e., using primers that only amplify the species of interest). By contrast, characterization of
9 microbial communities from environmental samples¹⁰ are often conducted using “universal” primers
10 amplifying bacterial 16S ribosomal RNA genes (rRNA) that yield, after DNA sequencing, enough
11 sequence information to identify the species carrying each DNA molecule.

12 Here, we use next-generation sequencing to characterize PCR products amplified from each eDNA
13 sample using taxon-specific primers to obtain a quick and cost-efficient assessment of the
14 macroorganism diversity. First, we describe PrimerTree, a novel R package that enables *in silico*
15 evaluation of the specificity and information content of “universal” primer pairs. Next, we describe the
16 application of this approach to the analysis of eDNA extracted from 91 samples of 40 mL of surface
17 water collected along the Cuyahoga River (Ohio, USA). We amplified each sample using 12 primer pairs
18 targeting mammals¹¹, fish¹², amphibians¹², birds¹³, bryophytes¹³, arthropods¹⁴, copepods¹⁵ and plants¹⁶
19 (as well as several microorganism taxa^{13,17,18}) and, after indexing, sequenced simultaneously the PCR
20 products by massively parallel sequencing. Our analyses show that this methodology can provide a
21 broad perspective on the aquatic and terrestrial biodiversity at the sampled sites in a simple, rapid and
22 cost-effective manner.

23

1 **Materials and Methods**

2 *In-silico evaluation of universal primer pairs*

3 To evaluate the amplification breadth and informativity of “universal” primers we developed
4 PrimerTree, an R package that performs the following functions for each primer pair provided by the
5 user:

- 6 (1) *In silico* PCR against a selected NCBI database
- 7 (2) Retrieval of DNA sequences predicted to be amplified
- 8 (3) Taxonomic identification of these sequences
- 9 (4) Multiple DNA sequence alignment
- 10 (5) Reconstruction of a phylogenetic tree
- 11 (6) Visualization of the tree with taxonomic annotation

12 PrimerTree utilizes the *in silico* primer search implemented in Primer-BLAST¹⁹ by directly querying the
13 NCBI Primer-BLAST search page. This allows access to all options available on the NCBI website. By
14 default, PrimerTree searches the NCBI (nt) nucleotide database but alternative NCBI databases, such as
15 only assembled genomes, or Refseq mRNA, can be queried. Note that when the proposed primers are
16 degenerate, PrimerTree automatically tests up to 25 possible combinations (by default, with more
17 possible) of primer sequences in Primer-BLAST and merges the results. The primer alignment results are
18 then processed using the NCBI E-utilities (<http://www.ncbi.nlm.nih.gov/books/NBK25500/>) to i) retrieve
19 DNA sequences located between the primers (i.e., the “amplified” sequences) and ii) obtain taxonomic
20 information related to each DNA sequence using the NCBI taxonomy database
21 (<http://www.ncbi.nlm.nih.gov/books/NBK21100/>). PrimerTree next aligns all “amplified” DNA
22 sequences using Clustal Omega²⁰ with a user configurable substitution matrix and reconstructs a
23 Neighbor-Joining tree using the ape package²¹. Finally, PrimerTree displays the resulting phylogenetic

1 tree using the `ggplot2` package, labeling each taxon in a different color and adding the names of the
2 main taxa using the `directlabels` package (<http://CRAN.R-project.org/package=directlabels>)²².
3 PrimerTree usually runs in less than five minutes, but the runtime varies greatly depending on the
4 primer specificity (i.e., how many DNA sequences are “amplified”), the search parameters chosen, the
5 current load on the NCBI servers and internet connection. In particular, each degenerate position in a
6 primer will result in up to four times as many primers to be tested, which can considerably increase the
7 runtime. To limit maximum runtime in this situation, PrimerTree randomly samples only a portion of the
8 total primer permutations (25 by default). Changing the number of sampled permutations or including
9 all variants is possible. PrimerTree uses the `plyr` package ([https://cran.r-](https://cran.r-project.org/web/packages/plyr/index.html)
10 [project.org/web/packages/plyr/index.html](https://cran.r-project.org/web/packages/plyr/index.html)) extensively and has full support for any of the parallel
11 backends compatible with the `foreach` package ([https://cran.r-](https://cran.r-project.org/web/packages/foreach/index.html)
12 [project.org/web/packages/foreach/index.html](https://cran.r-project.org/web/packages/foreach/index.html)). In particular, parallel retrieval of the primer sequences
13 from NCBI speeds up the total runtime considerably. Note that parallel queries to Primer-BLAST are
14 queued by NCBI’s servers and are only processed once there is free compute time.

15

16 *Sampling and DNA extraction*

17 We collected water samples from the upper (n=39), middle (n=16) and lower (n=24) Cuyahoga River
18 (**Figure 1** and **Supplemental Table 1** for details). These sections correspond to, respectively, an area with
19 lower population density; a section with more human presence, dams and water treatment plants; and
20 a heavily industrialized area. Each sample consisted of ~50 mL of surface water collected roughly one
21 meter from the bank of the river in a sterile 50ml conical tube. In addition, we also collected additional
22 samples from water sources entering the river (n=12, **Supplemental Table 1**). The collections were
23 performed during two sampling times separated by a major rain episode that dramatically increased the

1 river discharge (**Supplemental Figure 1**). Samples were stored at ambient temperature until return to
2 lab, where they were centrifuged and frozen at -20°C.
3 We isolated DNA from each water sample using the following procedure adapted from previous
4 studies^{23,24}. We first mixed each water sample by inversion and transferred 40 mL to a new tube for
5 centrifugation at 8,000 x *g* for 30 minutes at 4°C. We discarded the supernatant and resuspended the
6 pellet in 1 mL of ATL lysis buffer (DNeasy kit, Qiagen) supplemented with 0.47% Triton-X (Ricca Chemical
7 Company), 7.88 mg of lysozyme (Fisher Scientific) and 19.2 units of lysostaphin (Sigma Aldrich). We then
8 incubated the samples at 37°C for 1 hour while shaking them. We digested further by incubating 350 µl
9 of each sample with 50 µl of proteinase K and 350 µl buffer AL (Qiagen) for 60 minutes at 56°C. Finally,
10 we extracted DNA using Qiagen DNeasy columns according to the manufacturer's instructions. We
11 included two extraction controls and processed them identically and at the same time as the rest of the
12 samples to monitor cross- or laboratory contamination. In addition, all experiments were performed in a
13 laboratory where no eDNA or vertebrate DNA (aside from human and mouse) had been previously
14 extracted or amplified.

15

16 *DNA amplification and sequencing*

17 We selected primers to amplify different taxa from the literature and evaluated *in silico* their specificity
18 and information content using PrimerTree. We required each primer set to amplify a region small
19 enough for MiSeq reads to overlap to enable correction of sequencing errors. We amplified DNA
20 extracted from each sample (with two extraction controls and one PCR negative control) using the
21 following conditions: initial denaturation of 95°C for 15 minutes followed by 50 cycles of 95°C for 30
22 sec., 55°C for 30 sec. 72°C for 30 sec in 1X Quantitect mastermix (Qiagen) with 0.4 µM of each primer.
23 Overall, we performed, on each sample, 12 independent DNA amplifications targeting Archaea (16S

1 rRNA)¹⁷, Mammals (16S rRNA)¹¹, Algae (23S rRNA)²⁵, Amphibians (mt-Cytb)¹², Birds (12S rRNA)¹³, Fish
2 (mt-Cytb)¹², Bryophytes (trnL)¹³, Arthropods (mt-Co1)¹⁴, Copepods (28s rRNA)¹⁵, Diatoms (18S rRNA)¹⁸,
3 Fungi (ITS)¹³ and vascular plants (trnL)¹⁶ (**Supplemental Table 2**). Each primer included a 5' tail for
4 barcoding and Illumina sequencing (see below). We then pooled all 12 amplification products obtained
5 from each water sample. We added Illumina adapter sequences and labeled each sample with an
6 individual six nucleotide index (with each index distinct from all other index by at least two nucleotides)
7 using primers targeting the 5' oligonucleotide tail in 10 cycles of PCR (initial denaturation of 94°C for 3
8 minutes followed by 10 cycles of 94°C for 45 sec., 56°C for 45 sec. 72°C for 45 sec in 1X buffer, 1.25U
9 GoTaq (Promega), 2mM MgCl₂ and 2μM of each primer). We then pooled all indexed samples together
10 and we sequenced the resulting library on an Illumina MiSeq to generate 10,507,986 paired-end reads
11 of 250 bp.

12

13 *DNA sequence analysis pipeline*

14 We used custom PERL scripts to retrieve the index information and identify and trim the amplification
15 primer sequences. We discarded any DNA sequence shorter than 50 bp, with the exception of
16 sequences amplified with plant and bryophyte primers for which short amplification products were
17 expected (see **Supplemental Table 1**). We also discarded any trimmed read pair for which the difference
18 in sequence length was greater than 5 bp between the two paired-end reads to eliminate any reads
19 where primers were not found in both reads. We then merged the paired reads into a single consensus
20 DNA sequence using PANDAseq with default parameters²⁶. We used Mothur²⁷ to cluster unique DNA
21 sequences and counted how many reads carried each unique DNA sequence.

22 While all raw sequences are freely available online (accession number SRP058316), we only describe
23 here the analyses of macroorganism DNA sequences for sake of simplicity (microorganism sequences

1 can be easily analyzed using standard packages such as those implemented in QIIME²⁸). For macro-
2 organisms such as mammals, many species have been sequenced for the locus of interest and if not, a
3 closely related species is likely present in the NCBI database (but see also below). Therefore, to analyze
4 DNA sequences from macroorganisms – mammals, amphibians, birds, bryophytes, arthropods,
5 copepods and plants – we used BLAST²⁹ to directly identify the closest DNA sequences in the NCBI
6 database and the likely species of origin. Briefly, we removed from our analyses any DNA sequence
7 observed in less than 10 reads total (summing across all samples), as these likely represent sequencing
8 errors. We then compared each remaining DNA sequence to all sequences deposited in the NCBI nt
9 database using Blastn (excluding uncultured samples) and only considered matches with greater than
10 90% identity over the entire sequence length. We then retrieved taxonomic data of all best match(es)
11 for each sequence from NCBI. If multiple species matched a single sequence, all species names were
12 assigned to the sequence. We conducted further analyses at the species level for all taxa, using a
13 minimum read count per sample of 10 to determine absence/presence.

14

15 **Results**

16 *In silico assessment of universal primer pairs using PrimerTree*

17 PCR primers are usually designed to amplify one locus in a specific organism. Even “universal” primers,
18 designed to amplify many species within the same taxonomic group, are typically used to only amplify
19 DNA extracted from a single macroorganism (of this taxon). For studies amplifying eDNA or DNA from
20 unknown taxa, we require that the amplification works on all members of a given taxon while avoiding
21 off-target amplification (that could reduce the number of sequences from the desired taxon). In
22 addition, the amplified region should contain enough sequence information to identify the species
23 carrying the DNA sequence.

1 PrimerTree enables a rapid and visual assessment of these parameters for any primer pair by displaying
2 the results of *in silico* PCR as a taxonomically-annotated phylogenetic tree. For example, **Figure 2** shows
3 a subset of the PrimerTree results for a primer pair targeting the mammalian mitochondrial 16S
4 ribosomal RNA genes (**Figure 2A**) and primers designed to amplify the chloroplast trnL gene of non-
5 vascular plants (**Figure 2B**). The tree display enables rapid evaluation of the specificity of the primer
6 pairs (e.g., off-target amplification of amphibians and ray-finned fishes on **Figure 2A**). In addition, the
7 information content can be easily assessed by the length of the branches leading to different sequences
8 (scaled in number of nucleotide differences). For example, PrimerTree reveals much longer branch
9 lengths on **Figure 2A** than in **Figure 2B** suggesting a better discriminating power for the mammalian
10 sequences than the bryophyte sequences (see also below). By default, PrimerTree displays phylogenetic
11 trees annotated at all taxonomic levels enabling the user to determine the level of specificity of each
12 primer set (**Supplemental Figure 2**).

13 *High-throughput sequencing of DNA amplified from river samples*

14 We analyzed DNA extracted from 40 mL of surface water collected in 91 sites along the Cuyahoga River
15 (**Figure 1**). For each sample, we performed 12 PCR amplifications targeting mammals, amphibians, birds,
16 fish, arthropods, copepods, bryophytes and vascular plants (as well as several microorganism taxa not
17 analyzed here). We then individually indexed the PCR products of each sample and sequenced them on
18 an Illumina MiSeq (**Figure 3**) to generate a total of 10,507,986 paired-end reads (**Table 1**). After stringent
19 quality filtering we retained between 1,645,452 and 6,213 reads for the analysis of each taxon (**Table 1**).

20 *Species identification and resolution*

21 In contrast to microorganisms, where several hundred species are likely present in a given sample, we
22 only expect to amplify DNA from a few different macro-organism species per sample. Therefore, even if
23 the number of initial DNA molecules from a species is low, as long as the template can be amplified

1 efficiently, many reads will be generated from these few DNA sequences. For example, if a sample had
2 DNA from ten mammals and one species only accounted for 1% of the total DNA, its DNA sequence
3 would be represented, on average, by 142 reads (mammal amplifications were represented by, on
4 average, 14,211 reads per sample). We therefore considered that DNA sequences represented by less
5 than 10 reads total across all samples were caused by sequencing errors and discarded them, removing
6 between 1% and 25% of all reads generated, depending on the taxon considered (**Table 1**). We blasted
7 the remaining DNA sequences to identify the closest DNA sequences in NCBI and to assign a species
8 label to each DNA sequence.

9 In agreement with our *in silico* analyses, we observed large variations among primers in the specificity of
10 the taxa identified (**Supplemental Table 3**). For some primers, the sequence information was insufficient
11 to differentiate the organisms down to the species level: for example, each DNA sequence amplified
12 from the *trnL* gene of vascular plants matched sequences from 34.89 different species on average
13 (**Supplemental Table 4**). In fact, DNA sequences amplified from this primer pair matched a single taxon
14 only when considering families or higher taxonomic levels. This contrasted with more informative DNA
15 sequences such as the mammalian 16S rRNA for which each DNA sequence generated matched, on
16 average, only 1.27 species (**Supplemental Table 3**). Note that the observed specificity of the bird primers
17 differed from that expected from *in silico* analyses (12.39 species per DNA sequence generated
18 compared to 1.57 expected based on sequences available in NCBI). This apparent low specificity in our
19 data was caused by the presence of many DNA sequences from thrushes (a family of passerine birds) for
20 which many species with the exact same DNA sequences at the 12S rRNA gene have been sequenced. In
21 addition to differences in their information content (that influences the ability to identify a given
22 sequence), primer pairs also differed in the amplification specificity. For example, while initially designed
23 to amplify mammals¹¹, the mammalian 16S rRNA primers also amplified *Actinopterygii* (ray-finned

1 fishes). On the other hand, the primers targeting cytochrome B of fish¹² only amplified common carp
2 (*Cyprinus carpio*) and creek chub (*Semotilus atromaculatus*) under our PCR conditions as predicted by
3 our PrimerTree analysis (**Figure 2C,D and Supplemental Table 5**). The “mammalian” 16S rRNA primers
4 successfully amplified 31 samples for common carp while 15 were amplified using the cytochrome B
5 primers, indicating a higher sensitivity (under our PCR conditions). The samples positive for common
6 carp by the cytochrome B primers were usually also positive by the 16S rRNA primers: of the 15 samples
7 positive for common carp using cytochrome B primers, 11 were positive using 16S rRNA primers
8 (**Supplemental Table 6**). The four samples that did not yield common carp DNA using rRNA primers (but
9 were positive using cytochrome B) could perhaps be explained by very low amount of fish DNA
10 molecules (and stochastic amplification) or by the presence of much more abundant amplifiable DNA
11 (e.g. mammalian DNA) that might have entirely swamped this signal. Consequently, for all subsequent
12 analyses we used fish sequences amplified by the mammalian 16S RNA primers rather than those
13 amplified with the cytochrome B fish primers.

14 *Molecular assessment of the biodiversity of the Cuyahoga River*

15 Overall, across 91 water samples collected along the Cuyahoga River, we identified 54 samples positive
16 for fish DNA (representing 15 species), 77 samples positive for mammalian DNA (17 species excluding
17 human), 12 samples positives for bird DNA (from at least eight species), 18 samples positive for
18 arthropod DNA (15 species), 16 samples positive for copepod DNA (two species) while the “amphibian”
19 primers amplified turtle and two-lined salamander DNAs in two samples (**Table 2, Supplemental Table**
20 **6**). In addition to many organisms living in the river (e.g., fish, aquatic insects) or semi-aquatic animals
21 (beaver, mink, muskrat), we also amplified DNA from many terrestrial species that live near the banks of
22 the Cuyahoga River such as raccoon, groundhog, squirrel or mouse. Similarly, we identified DNA from
23 many birds that live on (swan, duck, sea gulls) as well as near the river (sparrow, wild turkey). The

1 species identified often corresponded with the local environment where the samples were collected: for
2 example, beaver DNA was amplified from samples collected in protected forested areas, gull DNA near
3 Lake Erie. In this regard, it is interesting to note that fish DNA showed significant differences in their
4 geographical distribution. DNA from fish of the *Moxostoma* genus (probably Silver Redhorse) was
5 commonly detected in the Upper Cuyahoga River but rare elsewhere ($p=0.02$, **Figure 1**). The central
6 stoneroller was detected only in the middle and lower Cuyahoga ($p=3.1 \times 10^{-3}$). On the other hand,
7 common carp were found throughout the entire Cuyahoga River ($p=0.23$). Surprisingly, we also
8 identified DNA from one invasive Asian Carp species in the Cuyahoga River near Lake Erie (**Figure 1**).
9 Note that the extraction and PCR controls almost exclusively yielded human DNA sequences (>99 % of
10 the reads) with one extraction control also displaying pig DNA in 0.6 % of the reads. It is important to
11 note here that, since the products of all PCRs were pooled and sequenced, regardless of the presence of
12 detectable amplified products on an agarose gel, the sensitivity of this approach is magnified compared
13 to standard molecular approaches and laboratory contamination with human DNA difficult to avoid.
14 However, these controls indicated that DNA sequences (aside from human sequences) retrieved from
15 the water samples must be genuine and that cross-contamination in the laboratory was minimal.

16

17 *Temporal variations in diversity*

18 We performed the collection of water samples at two time points separated by an episode of heavy rain
19 falls that dramatically altered the water level and flow of the Cuyahoga River: the discharge at Hiram
20 Rapids, in the Upper Cuyahoga River, increased from 2.46 cubic meters per second (close to the median
21 daily statistics) at the time of the first sampling to 20.87 cubic meters per second and was still as high as
22 high as 5.95 cubic meters per second at the time of the second sampling 12 days later (**Supplemental**
23 **Figure 1**). For most macro-organisms, we did not detect any statistical difference between the samples

1 collected at the two time points. One notable exception concerns fish DNA: we observed significantly
2 more positive samples before rain than after ($p=0.05$, Fisher's exact test).

3

4 **Discussion**

5 Since the first reports that DNA could be retrieved from environmental samples³⁰, studies of
6 environmental DNA have broadened in scope from studies of bacterial communities (e.g.,^{2,31}), to
7 identification and monitoring of a given species (e.g.,^{5,9,12,32}) and the characterization of microorganism
8 populations (e.g.,^{33,34}) and recently to the identification of macroinvertebrates or vertebrates (e.g.,³⁵⁻³⁸).
9 However, several factors have limited a broader implementation of these approaches for ecological
10 studies. These limitations include the lack of tools enabling a wide range of species to be studied
11 simultaneously, the amount of starting material (often liters of water for aquatic environments) and the
12 high costs of such analyses. Additionally, the careful evaluation of potential primer pairs prior to
13 laboratory work is critical as most published primers have only been tested on DNA directly extracted
14 from the target organism and, while efficient, might not be specific or could even better amplify other
15 organisms.

16 *PrimerTree provides a robust assessment of primer pairs for eDNA studies*

17 We present here a simple R package that enables rapid screening of primers suitable for differentiating
18 species from a chosen taxon. PrimerTree allows evaluating the specificity of a given primer pair (to avoid
19 off-target amplification) and whether the amplified DNA sequences would provide enough information
20 to identify the organisms carrying the DNA sequences. Since the amplified regions are typically short
21 (100-300 bp), the resulting phylogenetic trees (**Figure 2** and **Supplemental Figure 2**) do not necessarily
22 reflect the true species relationships, but they enable an easy and rapid assessment of the primer
23 suitability. First, the automatic taxonomic annotation of the branches allows the user to quickly see

1 which taxa present in the BLAST database are amplifiable by the primers. This enables identifying that
2 undesirable (off-target) taxa may be amplified or that the primers might not amplify specific genera
3 within the taxon of interest. Second, the branch lengths show how many nucleotides separate DNA
4 sequences from different species. Short branch lengths might results in a low resolution in taxon
5 identification (or elevated risk of false-positives caused by sequencing or PCR errors) while long
6 branches will lead to high confidence species identification (e.g., **Figure 2A** vs. **2B**).

7 An attractive feature of PrimerTree is its ease of use. As an R package, it is easily installed or updated
8 and only requires that clustal be installed in the user's path. The entire process of getting BLAST hits for
9 a primer pair is done through R using a single command and plots are generated using a second
10 command. The summary function provides many useful statistics including amplified DNA sequence
11 lengths, number of taxa amplified and average pairwise differences within each taxonomic level.
12 Additionally, information on the BLAST results, sequences obtained, taxonomy of amplifiable sequences,
13 neighbor joining distance matrix and phylogenetic tree are preserved within the PrimerTree object if the
14 user wants to further analyze or further summarize the results. Finally, PrimerTree runs fast: for a
15 primer pair with no degenerate bases, PrimerTree results are usually obtained in less than five minutes.

16 *Limitations of PrimerTree and comparison with other programs*

17 One limitation of PrimerTree is that, by default, it only retrieves up to 500 random amplifiable
18 sequences from the NCBI database and this random subset might not accurately represent the
19 specificity of a given primer pair for a particular environment. For example, the bird primers used in our
20 study display a high species-specificity for most birds but were not able to differentiate among thrushes.
21 Note however that this issue can be easily circumvented by selecting specific targets in PrimerTree if
22 one is interested in a particular taxon (e.g., one could run PrimerTree for querying only *Turdidae*
23 sequences in NCBI to evaluate if a given primer pair is informative within this taxon). Additionally, the

1 problem of returning a random subset of all sequences can be minimized by requesting more
2 amplifiable sequences using the `num_aligns` argument in the command. Another important limitation is
3 that PrimerTree does not highlight taxonomic groups that are not retrieved. If a particular taxonomic
4 group is not amplifiable by a primer pair, the user must recognize the absence of that taxon. For
5 example, if one primer pair amplifies all mammals except monotremes, the user must note the absence
6 of the order *Monotremata* from the resulting plot or taxonomic information in the PrimerTree object.
7 Along these same lines, the absence of a species or larger taxonomic group from PrimerTree results may
8 not mean that the assayed primers cannot amplify those taxa. As PrimerTree uses the BLAST nr/nt
9 database as a reference, any species without DNA sequence for the targeted locus in the database
10 cannot be retrieved.

11 EcoPCR³⁹ is another program than can evaluate primer specificity by utilizing a user-provided database
12 to identify amplifiable sequences and provides a summary of amplifiable sequences and their taxonomic
13 information. The ecoPCR program retrieves all amplifiable sequences in a custom database, compared
14 to only a subset of the BLAST nr/nt database for PrimerTree. However, one advantage of PrimerTree is
15 that it always queries the most current available version of the entire BLAST nr/nt database, while the
16 user would need to update the database provided to ecoPCR and re-generate the ecoPCR database files
17 to include new sequences in the analysis. Additionally PrimerTree benefits from the graphical output,
18 which is automatically generated using the `plot` function in R, to summarize the results in a user-friendly
19 manner.

20 *Characterization of eDNA extracted from the Cuyahoga River*

21 We analyzed eDNA extracted from 91 samples of water collected along the Cuyahoga River. We showed
22 that, despite a small amount of starting material, we were able to retrieve eDNA from a wide range of
23 organisms, including many vertebrate species: on average, each water sample contained eDNA from 3.2

1 vertebrates (fish, mammal or bird), 0.3 arthropods and numerous plants. The most exciting feature of
2 our findings is the amplification of DNA from many terrestrial organisms from the river samples. We
3 amplified DNA sequences from deer, squirrel, raccoon, groundhog, vole, mink and skunk in addition to
4 semi-aquatic species such as beaver and muskrat. We also identified DNA sequences for agricultural
5 species (cow and swine) and companion animals (dog and cat), consistent with the findings of a study
6 recently made available through bioRxiv (<http://dx.doi.org/10.1101/020800>). DNA from these species in
7 the river samples likely originated from fecal matter washing into the river, deceased animals in or
8 around the river or animals drinking or wading in the river. Similarly, we detected many avian
9 sequences, including species that do not live in the river (e.g., chickens, turkeys, blackbirds and
10 sparrows).

11 Several species were unevenly distributed along the Cuyahoga, such as the Silver Redhorse that was
12 mostly detected in the Upper Cuyahoga, which demonstrates the potential of our approach to identify
13 variations in species distribution (**Figure 1**). Among the 48 water samples positive for fish DNA, we also
14 identified one sample containing DNA from one Asian carp (either *Ctenopharyngodon idella* or
15 *Mylopharyngodon piceus*), invasive species of the Great Lakes that were not known to be present in the
16 Cuyahoga River (**Figure 1**).

17 *Biological and technical factors influencing species detection and the rate of false negatives*

18 Our two samplings of the Cuyahoga River were separated by a major rain event. Since we sampled the
19 upper portion of the Cuyahoga both before and after the rain, we compared species detection rates to
20 determine if the influx of water into the river would stir up sediment and wash genetic material into the
21 river to increase species detection rates or dilute the river water and decrease species detection rates.
22 We found that the number of samples positive for mammals, birds, arthropods and copepods did not
23 change after the rain, whereas the number of samples positive for fish species decreased after the rain.

1 This may suggest that, for fish, the rainwater dilutes the genetic material of fish, but for other taxa this
2 dilution effect might be balanced by the influx of genetic material in water washing into the river.
3 Replication of these results using other sample sets is necessary to confirm that this is a general
4 phenomenon.

5 Certain taxa were identified less frequently than we would have expected (e.g., amphibians, arthropods
6 and copepods, **Table 1**). This could be due to a number of factors. Our DNA extraction protocol (relying
7 on the DNeasy extraction kit after centrifugation) may not have been optimal to retrieve DNA from such
8 organisms as the choice of DNA isolation protocol from water samples can strongly influence the
9 proportion of DNA isolated from different taxa⁴⁰. In particular, it is important to note that our isolation
10 method was unlikely to capture free DNA in the river water sample, but rather only DNA within, or
11 adhered to, particulate matter which could lead to disproportional representation of the eDNA (see e.g.,
12 ^{24,41}. It is also possible that the PCR conditions we used were inadequate to efficiently amplify some of
13 the targeted templates^{42,43}. Alternatively, the amount of genetic material present in the river may be
14 lower for these taxa than for others such as fish or mammals (possibly due to differences in the size of
15 the dead animals or the amount of feces or shed tissues). Additional studies, including experimental
16 validation of the primers for these taxa, will be necessary to differentiate these possibilities.

17 While detection of eDNA is by nature stochastic, the small sample size used in our study likely increases
18 this randomness and increases the chance of false negatives. By optimizing PCR conditions and
19 increasing the number of cycles, one can limit the chance of having a DNA template present in a given
20 sample failing to amplify. Nonetheless, a small sample of water may not entirely capture the diversity
21 present in the environment, especially for species represented by few DNA molecules in the
22 environment. While the amplification of a species' DNA sequence is a clear evidence of the presence of
23 this species (assuming a low level of cross-contamination), the failure to detect an organism is not a

1 proof of its absence in the environment as the false negative rate is likely to be important. For example,
2 while 20 out of 38 samples in the upper portion contained common carp (*Cyprinus carpio*) DNA, we
3 would not be able to rule out the presence of common carp at the other sampled locations in this part
4 of the river. The use of universal primer pairs also contributes to the likely high false negative rates:
5 since our assay is based on primers that amplify multiple taxa rather than being species-specific, it is
6 possible that DNA sequences from one species completely overwhelmed the signal originating from
7 rarer eDNA templates present in the same sample. This effect could be magnified if there are
8 differences in amplification efficiencies^{42,43}. For example, it is possible that some samples that did not
9 yield common carp sequences actually contained common carp eDNA but that these molecules
10 remained undetected due to the abundance of other templates (e.g., Silver Redhorse or human) in these
11 samples. The presence of multiple amplification targets is also one reason why this assay is non-
12 quantitative: the number of reads obtained for one species in a given sample is not only determined by
13 how much DNA from that species is present in the sample, but also by how many other amplifiable
14 species are present. Therefore, one cannot directly compare species read counts between samples to
15 determine relative abundance³⁶. In these regards, it is important to emphasize that the present
16 approach does not replace classical eDNA studies targeting a single species but is designed for enabling
17 broader ecological survey (that may guide further in-depth investigations).

18 *A fast, high-throughput and cost-efficient method to characterize biodiversity from environmental*
19 *samples*

20 We described here a customizable approach that builds on existing methods⁴²⁻⁴⁵ to enable simultaneous
21 analyses of eDNA extracted from many samples for a wide range of taxa.

22 We showed that this methodology can be successfully applied to analyze small water samples. This
23 aspect, while it has its limitations (see above), is essential for any study for which sampling and storing

1 of large volumes is logistically challenging or impossible. In particular, small volume collection enables
2 sampling in sites that can only be reached by hiking or paddling and will therefore make this approach
3 suitable for many ecological studies. Additionally, sample processing (such as filtration) in the field can
4 be slow and requires specialized equipment, particularly for large sample volumes. Here, we chose
5 instead to centrifuge the samples upon return to the laboratory. This approach isolates the same
6 portion of the sample as filtration (particulate matter), is rapid and reduces potential sample
7 contamination during filtration. Additionally, by minimizing in-field handling, more samples can be
8 collected across an aquatic system in a given timeframe, minimizing temporal sampling bias. For
9 instance, we were able to sample roughly half the length of the Cuyahoga River in a single day with a
10 single sampling team using kayaks. The speed of sampling is also an important parameter as it enables
11 studies investigating the consequences of specific events such as rainfall, ecological disasters (e.g., oil,
12 chemical or pollution spills) or even transient ecological events such as fish spawning. However,
13 depending on the specific goals of a study and the resources available, larger sample volumes or
14 alternate DNA isolation methods might be preferable.

15 Finally, this approach is cost-efficient: in our study of the Cuyahoga River, we characterized eDNA
16 amplified from mammals, fish, birds, amphibians, arthropods, vascular plants, bryophytes and many
17 microorganisms in 91 water samples for a total cost of approximately US \$2000 (including DNA
18 extraction, PCR reagents and sequencing costs, see **Supplemental Table 7** for details). This high level of
19 multiplexing (across samples and across taxa) decreases the burden of the next-generation sequencing
20 price, but, thanks to the tremendous sequencing output, still provides enough reads to rigorously
21 characterize the composition of each sample. The use of an inexpensive PCR to add the sequencing and
22 barcoding adapters also dramatically reduces the cost compared to generation of typical next-

1 generation sequencing libraries^{18,32,36,46}. Note that the cost and efficiency could potentially be further
2 improved by multiplexing primer pairs in the PCR reaction⁴⁴.

3 *Limitations and cautionary notes for future studies*

4 One important limitation of our approach regards the interpretation of the results. First, because we
5 performed the species identification of sequences using BLAST, it is affected by the content of the NCBI
6 database and the reliability might vary for different taxa. For instance, our analyses revealed, in one
7 sample, the presence of DNA sequences most similar to a Taiwanese vole, which is not present near the
8 Cuyahoga River, but with only 93.4% identity. These sequences likely originate from a local species of
9 vole (closely related to the Taiwanese vole) that has not been sequenced for the 16S rRNA gene. This
10 illustrates that matches with low percent identity BLAST hits need to be cautiously interpreted. Using
11 primer pairs targeting different loci (e.g., COI) could partially circumvent this limitation as one species
12 may have been sequenced at one locus but not another. On the other hand, this example also shows
13 that this method can be used to identify organisms even if they have not been previously sequenced for
14 the locus of interest, as long as a closely related species is present in the database, but that rigorously
15 identifying the actual species present will require further analyses. In addition, when considering best
16 match for a given sequence, the identification should be considered in the context of the other
17 sequences amplified from the sample. For instance, in our analyses we identified sequences that best
18 matched unexpected fowl species (e.g., *Gallus lafayetii* and *Gallus sonneratii*). However, these DNA
19 sequences were only observed in small numbers and in samples that had a much larger numbers of
20 *Gallus gallus* sequences. The most likely explanation for these sequences is that these reads derived
21 from PCR or sequencing errors from *Gallus gallus* templates: the chicken sequences were so abundant in
22 these samples that even rare errors led to a number of reads sufficient to pass through our stringent

1 filtering criteria. The same phenomenon generated a handful of sequences most similar to monkey or
2 ape sequences in samples with many human reads.

3 One final issue is the risk of false positives generated by contamination: since PCR products are directly
4 sequenced by massively parallel sequencing, even minute contamination can be detected and analyzed.

5 It is therefore critical to implement stringent measures to prevent and detect contamination, similar to
6 those used for ancient DNA studies (e.g., numerous negative controls, use of sterilized room and
7 equipment, etc), to prevent cross-contamination and obtain reliable results.

8

9 *Conclusion*

10 The study of environmental DNA using massively-parallel sequencing technologies enables
11 characterization of species biodiversity in a simple, high-throughput and cost-effective manner. Our
12 results reveal that a sampling protocol relying on small sample volume enables the preliminary
13 evaluation of complex environments. The ability to characterize a very diverse range of taxa from a
14 single sample in a high-throughput manner will allow future studies to expand the scope of the
15 biodiversity studied and to explore complex ecological interactions among species. Additionally, the
16 identification of eDNA of local terrestrial flora and fauna in river samples provides a simple way to assess
17 the local diversity of environments adjacent to rivers or other water bodies. Overall, our findings
18 illustrate the sensitivity and utility of broad surveys of eDNA by deep sequencing and shows how this
19 approach can constitute an excellent foundation for ecological and environmental studies.

20

1 **References**

- 2 1 Taberlet, P., Coissac, E., Hajibabaei, M. & Rieseberg, L. H. Environmental DNA. *Mol. Ecol.* **21**,
3 1789-1793, (2012).
- 4 2 Rusch, D. B. *et al.* The Sorcerer II Global Ocean Sampling expedition: northwest Atlantic through
5 eastern tropical Pacific. *PLOS Biol.* **5**, e77, (2007).
- 6 3 Venter, J. C. *et al.* Environmental genome shotgun sequencing of the Sargasso Sea. *Science* **304**,
7 66-74, (2004).
- 8 4 Minamoto, T., Yamanaka, H., Takahara, T., Honjo, M. N. & Kawabata, Z. i. Surveillance of fish
9 species composition using environmental DNA. *Limnology* **13**, 193-197, (2011).
- 10 5 Goldberg, C. S., Pilliod, D. S., Arkle, R. S. & Waits, L. P. Molecular detection of vertebrates in
11 stream water: a demonstration using Rocky Mountain tailed frogs and Idaho giant salamanders.
12 *PLOS One* **6**, e22746, (2011).
- 13 6 Biggs, J. *et al.* Using eDNA to develop a national citizen science-based monitoring programme for
14 the great crested newt (*Triturus cristatus*). *Biol. Conserv.* **183**, 19-28, (2015).
- 15 7 Pilliod, D. S., Goldberg, C. S., Arkle, R. S. & Waits, L. P. Factors influencing detection of eDNA
16 from a stream-dwelling amphibian. *Mol. Ecol. Resour.* **14**, 109-116, (2014).
- 17 8 Tréguier, A. *et al.* Environmental DNA surveillance for invertebrate species: advantages and
18 technical limitations to detect invasive crayfish *Procambarus clarkii* in freshwater ponds. *J. Appl.*
19 *Ecol.* **51**, 871-879, (2014).
- 20 9 Mahon, A. R. *et al.* Validation of eDNA surveillance sensitivity for detection of Asian carps in
21 controlled and field experiments. *PLOS One* **8**, e58316, (2013).
- 22 10 Kuczynski, J. *et al.* Microbial community resemblance methods differ in their ability to detect
23 biologically relevant patterns. *Nat. Meth.* **7**, 813-819, (2010).

- 1 11 Taylor, P. G. Reproducibility of ancient DNA sequences from extinct Pleistocene fauna. *Mol. Biol. Evol.* **13**, 283-285, (1996).
- 2
- 3 12 Thomsen, P. F. *et al.* Monitoring endangered freshwater biodiversity using environmental DNA. *Mol. Ecol.* **21**, 2565-2573, (2012).
- 4
- 5 13 Epp, L. S. *et al.* New environmental metabarcodes for analysing soil DNA: potential for studying past and present ecosystems. *Mol. Ecol.* **21**, 1821-1833, (2012).
- 6
- 7 14 Zeale, M. R., Butlin, R. K., Barker, G. L., Lees, D. C. & Jones, G. Taxon-specific PCR for DNA barcoding arthropod prey in bat faeces. *Mol. Ecol. Resour.* **11**, 236-244, (2011).
- 8
- 9 15 Bissett, A., Gibson, J. A. E., Jarman, S. N., Swadling, K. M. & Cromer, L. Isolation, amplification, and identification of ancient copepod DNA from lake sediments. *Limnol. Oceanogr.-Meth.* **3**, 533-542, (2005).
- 10
- 11
- 12 16 Taberlet, P. *et al.* Power and limitations of the chloroplast trnL (UAA) intron for plant DNA barcoding. *Nucleic Acids Res.* **35**, e14, (2007).
- 13
- 14 17 Baker, G. C., Smith, J. J. & Cowan, D. A. Review and re-analysis of domain-specific 16S primers. *J. Microbiol. Meth.* **55**, 541-555, (2003).
- 15
- 16 18 Zimmermann, J., Jahn, R. & Gemeinholzer, B. Barcoding diatoms: evaluation of the V4 subregion on the 18S rRNA gene, including new primers and protocols. *Org. Divers. Evol.* **11**, 173-192, (2011).
- 17
- 18
- 19 19 Ye, J. *et al.* Primer-BLAST: a tool to design target-specific primers for polymerase chain reaction. *BMC bioinformatics* **13**, 134, (2012).
- 20
- 21 20 Sievers, F. *et al.* Fast, scalable generation of high-quality protein multiple sequence alignments using Clustal Omega. *Mol. Syst. Biol.* **7**, 539, (2011).
- 22

- 1 21 Paradis, E., Claude, J. & Strimmer, K. APE: Analyses of Phylogenetics and Evolution in R language.
2 *Bioinformatics* **20**, 289-290, (2004).
- 3 22 Wickham, H. ggplot2: Elegant Graphics for Data Analysis. *Ggplot2: Elegant Graphics for Data*
4 *Analysis*, 1-212, (2009).
- 5 23 Chan, E. R. *et al.* A novel method for determining microflora composition using dynamic
6 phylogenetic analysis of 16S ribosomal RNA deep sequencing data. *Genomics* **98**, 253-259,
7 (2011).
- 8 24 Deiner, K., Walser, J.-C., Mächler, E. & Altermatt, F. Choice of capture and extraction methods
9 affect detection of freshwater biodiversity from environmental DNA. *Biol. Conserv.* **183**, 53-63,
10 (2015).
- 11 25 Sherwood, A. R. & Presting, G. G. UNIVERSAL PRIMERS AMPLIFY A 23S rDNA PLASTID MARKER IN
12 EUKARYOTIC ALGAE AND CYANOBACTERIA. *J. Phycol.* **43**, 605-608, (2007).
- 13 26 Masella, A. P., Bartram, A. K., Truszkowski, J. M., Brown, D. G. & Neufeld, J. D. PANDAseq:
14 paired-end assembler for illumina sequences. *BMC Bioinformatics* **13**, 31, (2012).
- 15 27 Schloss, P. D. *et al.* Introducing mothur: open-source, platform-independent, community-
16 supported software for describing and comparing microbial communities. *Appl. Environ.*
17 *Microbiol.* **75**, 7537-7541, (2009).
- 18 28 Caporaso, J. G. *et al.* QIIME allows analysis of high-throughput community sequencing data. *Nat.*
19 *Meth.* **7**, 335-336, (2010).
- 20 29 Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search
21 tool. *J. Mol. Biol.* **215**, 403-410, (1990).
- 22 30 Ogram, A., Sayler, G. S. & Barkay, T. The extraction and purification of microbial DNA from
23 sediments. *J. Microbiol. Meth.* **7**, 57-66, (1987).

- 1 31 Tringe, S. G. *et al.* Comparative metagenomics of microbial communities. *Science* **308**, 554-557,
2 (2005).
- 3 32 Ficetola, G. F., Miaud, C., Pompanon, F. & Taberlet, P. Species detection using environmental
4 DNA from water samples. *Biol. Lett.* **4**, 423-425, (2008).
- 5 33 Bik, H. M. *et al.* Metagenetic community analysis of microbial eukaryotes illuminates
6 biogeographic patterns in deep-sea and shallow water sediments. *Mol. Ecol.* **21**, 1048-1059,
7 (2012).
- 8 34 Hirai, J., Kuriyama, M., Ichikawa, T., Hidaka, K. & Tsuda, A. A metagenetic approach for revealing
9 community structure of marine planktonic copepods. *Mol. Ecol. Resour.* **15**, 68-80, (2015).
- 10 35 Andersen, K. *et al.* Meta-barcoding of 'dirt' DNA from soil reflects vertebrate biodiversity. *Mol.*
11 *Ecol.* **21**, 1966-1979, (2012).
- 12 36 Kelly, R. P., Port, J. A., Yamahara, K. M. & Crowder, L. B. Using environmental DNA to census
13 marine fishes in a large mesocosm. *PLOS One* **9**, e86175, (2014).
- 14 37 Mächler, E., Deiner, K., Steinmann, P. & Altermatt, F. Utility of Environmental DNA for
15 Monitoring Rare and Indicator Macroinvertebrate Species. *Freshwater Science* **33**, 1174-1183,
16 (2014).
- 17 38 Deiner, K. & Altermatt, F. Transport distance of invertebrate environmental DNA in a natural
18 river. *PLOS One* **9**, e88786, (2014).
- 19 39 Ficetola, G. F. *et al.* An in silico approach for the evaluation of DNA barcodes. *BMC Genomics* **11**,
20 434, (2010).
- 21 40 Rees, H. C. *et al.* REVIEW: The detection of aquatic animal species using environmental DNA - a
22 review of eDNA as a survey tool in ecology. *J. Appl. Ecol.* **51**, 1450-1459, (2014).

- 1 41 Turner, C. R. *et al.* Particle size distribution and optimal capture of aqueous microbial eDNA.
2 *Meth. Ecol. Evol.* **5**, 676-684, (2014).
- 3 42 Pinol, J., Mir, G., Gomez-Polo, P. & Agusti, N. Universal and blocking primer mismatches limit the
4 use of high-throughput DNA sequencing for the quantitative metabarcoding of arthropods. *Mol.*
5 *Ecol. Resour.* **15**, 819-830, (2015).
- 6 43 Elbrecht, V. & Leese, F. Can DNA-Based Ecosystem Assessments Quantify Species Abundance?
7 Testing Primer Bias and Biomass--Sequence Relationships with an Innovative Metabarcoding
8 Protocol. *PLOS One* **10**, e0130324, (2015).
- 9 44 Hajibabaei, M., Spall, J. L., Shokralla, S. & van Konynenburg, S. Assessing biodiversity of a
10 freshwater benthic macroinvertebrate community through non-destructive environmental
11 barcoding of DNA from preservative ethanol. *BMC Ecol.* **12**, 28, (2012).
- 12 45 Gibson, J. *et al.* Simultaneous assessment of the macrobiome and microbiome in a bulk sample
13 of tropical arthropods through DNA metasystematics. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 8007-
14 8012, (2014).
- 15 46 Cowart, D. A. *et al.* Metabarcoding is powerful yet still blind: a comparative analysis of
16 morphological and molecular surveys of seagrass communities. *PLOS One* **10**, e0117562, (2015).

17

18

19

1 Acknowledgments

2 This work was supported by Cleveland Clinic funds to DS.

3 Author Contributions

4 MVC, STS and DS conceived the study. MVC, JH, AS, ERC, KL, STS and DS contributed to the sample
5 collection. AS performed the DNA extraction and PCR amplification. MVC analyzed the data. JH
6 implemented PrimerTree. MVC and DS drafted the manuscript. All authors have read and approved the
7 final version.

8 Competing financial interests

9 The authors have no competing financial interests to disclose.

10 Data Accessibility

11 All sequencing data generated in this study are available through the SRA at NCBI (Accession
12 #SRP058316 (data to be released upon publication)). The R package PrimerTree is freely available at
13 <http://github.com/jimhester/primerTree>.

14
15

16 Tables and figures

17
18

19 Table 1: Sequencing and sequence analysis summary

20

Taxon targeted	Locus	Reads generated	Reads after QC	Unique seq.	Unique seq. w/ >= 10 reads	Reads represented by unique seq. >= 10	Species or OTUs identified	Seq. length (min-max)
Mammals	mt16S rRNA	1,384,149	1,352,524	13,238	786	1,328,072 (98.2%)	65	95 bp (76 - 118)
Amphibian	mt-Cytb	6,517	6,213	225	20	5,724 (92.1%)	2	68 bp (68 - 68)
Birds	12S rRNA	79,549	64,144	873	136	62,433 (97.3%)	123	53 bp (50 - 211)
Arthropods	COI	299,714	42,987	5,043	238	34,855 (81.1%)	25	162 bp (156 - 166)

Copepods	28S rRNA	502,972	79,486	13,252	368	59,753 (75.2%)	5	207 bp (53 - 368)
Fish	mt-Cytb	280,184	83,123	1,918	231	80,393 (96.7%)	5	90 bp (76 - 192)
Bryophytes	TrnL	253,053	225,780	1,728	262	223,153 (98.8%)	32	52 bp (25 - 55)
Vascular plants	TrnL	1,954,381	1,610,089	33,584	3,468	1,554,086 (96.5%)	236	48 bp (17 - 80)

1
2
3
4

Table 2: Macro-organisms identified in the Cuyahoga River.

Common name	Scientific name	Percent Identity	Samples positive
Mammals			
Cow	<i>Bos taurus/Bos indicus/Bos primigenius/Bos javanicus</i>	100	22
Dog	<i>Canis lupus/Canis aureus</i>	100	27
Deer	<i>Odocoileus virginianus/Odocoileus hemionus/Mazama americana</i>	100	2
Sheep	<i>Ovis aries/Ovis canadensis/Ovis vignei/Ovis dalli</i>	98.91	1
Beaver	<i>Castor canadensis</i>	100	16
Cat	<i>Felis catus/Felis silvestris</i>	100	3
Groundhog	<i>Marmota himalayana/Marmota monax</i>	97.78	2
Skunk	<i>Mephitis mephitis</i>	98.89	1
Mouse	<i>Mus musculus</i>	100	37
Mink	<i>Neovison vison</i>	100	4
Muskrat	<i>Ondatra zibethicus</i>	100	19
Raccoon	<i>Procyon lotor</i>	100	6
Rat	<i>Rattus norvegicus</i>	100	2
Squirrel	<i>Sciurus carolinensis</i>	97.83	6
Pig	<i>Sus scrofa/Sus barbatus/Sus philippensis</i>	100	26
Chipmunk	<i>Tamias striatus/Tamias sonomae</i>	100	4
Taiwan vole	<i>Microtus kikuchii</i>	93.41	1
Reptiles and amphibians			
Northern two-lined salamander	<i>Eurycea bislineata</i>	94.12	1
Box turtle	<i>Terrapene carolina</i>	100	1
Birds			
Teal	<i>Anas poecilorhyncha/Anas crecca/Anas platyrhynchos/Anas acuta/Anas clypeata/Mareca falcata/Cygnus melancoryphus/Histrionicus histrionicus/Ptaiochen pau/Thambetochen chauliodous/Anas hottentota/Anas querquedula/Anas cyanoptera</i>	100	4
Swan/Goose	<i>Cygnus columbianus/Anser fabalis/Anser anser/Cygnus olor/Branta bernicla/Callonetta leucophrys/Cygnus atratus/Anser indicus/Anser</i>	100	1

	albifrons/Anser cygnoides/Cygnus cygnus/Branta sandvicensis/Anser rossii/Branta canadensis/Anser canagica/Cygnus buccinator		
Chicken	Gallus gallus/Gallus sonneratii	100	7
Sparrow	Passer domesticus/Passer montanus/Chlorospingus canigularis	100	1
Ring-billed gull	Larus delawarensis	100	1
Wild turkey	Meleagris gallopavo	100	1
Great tit	Parus major	100	1
Thrush	Turdus sp.	100	1
Fish			
Central stoneroller	Campostoma anomalum	96.55	11
White sucker	Catostomus macrocheilus/Catostomus commersonii	100	13
Grass/Black carp	Ctenopharyngodon idella/Mylopharyngodon piceus	100	1
<i>Cyprinella spiloptera</i>	Cyprinella spiloptera	100	9
Common carp	Cyprinus carpio/Carassius gibelio/Cyprinus multitaeniata/Cyprinus melanes	100	35
Cypress minnow	Hybognathus hayi	94.78	1
Northern hogsucker	Hypentelium nigricans	100	10
Shiner	Luxilus chrysocephalus/Luxilus cornutus	100	5
Spotted sucker	Minytrema melanops	100	4
Redhorse	Moxostoma cervinum/Moxostoma anisurum	95.65	19
Bluehead chub	Nocomis leptocephalus	95.65	1
Emerald shiner	Notropis atherinoides	99.13	2
Blacknose dace	Rhinichthys atratulus	100	4
Bluntnose minnow	Pimephales notatus	100	9
Arthropods			
Non-biting midge/Fly	aff. Cyrtona/Tachinidae gen./Orthocladus sp./Orthoclaadiinae sp./Fannia serena	90.45	1
Crane fly	Antocha sp.	100	2
Fruit fly	Capparimyia aenigma/Scaptomyza frustulifera/Scaptomyza remota	90.45	3
Non-biting midge	Cricotopus bicinctus	100	1
Fruit fly	Drosophila medioimpressa	91.72	1
Whirligig beetle	Gyrinidae sp.	95.54	3
Non-biting midge	Microtendipes pedellus	99.36	2
Allegheny crayfish	Orconectes obscurus	100	1
Non-biting midge	Polypedilum convictum	100	1
Daphnia	Simocephalus cf.	99.36	3
Black fly	Simulium luggeri	100	1
Horsefly	Tabanus sp./Hybomitra zonalis	93.63	1
Crane fly	Tipula paludosa	96.82	1
Non-biting midge	Tribelos sp.	100	2
Non-biting midge	Xenochironomus xenolabis	98.09	1
Copepods			

1	Hemidiaptomus maroccanus/Hemidiaptomus amblyodon/Hemidiaptomus	94.1	16
	ingens/Hemidiaptomus roubau		
	Macrocyclops distinctus	93.73	2

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23

Figure 1. Geographic locations of water samples positive for Silver Redhorse and the Asian Carp eDNA using the 16S mammal primers. Each circle shows the location of a sampled site. The red circle indicates the location of the sample positive for Asian carp (*Ctenopharyngodon idella* or *Mylopharyngodon piceus*). Yellow circles are samples positive for Silver Redhorse (*Moxostoma cervinum* or *Moxostoma anisurum*). Map image was prepared by the Cleveland Clinic Center for Medical Art and Photography using Adobe Illustrator CS6 and points were overlaid using ggplot2.

Figure 2. Example of PrimerTree results. The figure shows phylogenetic trees annotated at the class level for (A) the mammalian 16S rRNA and (B) bryophyte trnL primer pairs. The complete PrimerTree results for all 12 primers used in this study are presented in Supplemental Figure 2. Panels C and D are PrimerTree results from mammal 16S rRNA and fish mt-Cytb primers, respectively, demonstrating a greater breadth of amplifiable sequences by the mammal 16S rRNA primers. 1,000 BLAST hits within *Actinopterygii* are shown and tips are colored by taxonomic order. Fish mt-Cytb PrimerTree results include only 6 orders, while mammal 16S rRNA primers can amplify 27. The order *Cypriniformes* (which includes carp) is green in both panels.

Figure 3. Experimental workflow. We first isolated DNA from 40 ml of river water. We then amplified each sample with 12 taxon-specific primer sets and pooled a portion of each PCR for each sample. Each primer had a 5' tail to allow a second PCR which added Illumina adapter sequence and an individual index. We then pooled all barcoded samples and sequenced the library on a MiSeq. We used the sequence information to identify species of origin for DNA fragments isolated from the original samples.



