

# NEON USER GUIDE TO MACROINVERTEBRATE DNA BARCODE (NEON.DP1.20126) AND ZOOPLANKTON DNA BARCODE (NEON.DP1.20221)

<b>PREPARED BY</b>	<b>ORGANIZATION</b>	<b>DATE</b>
Stephanie Parker	AQU	03/22/2018
Lee Stanish	FSU	03/22/2018

<i>Title:</i> NEON User Guide to Macroinvertebrate DNA Barcode (NEON.DP1.20126) and Zooplankton DNA Barcode (NEON.DP1.20221)	<i>Date:</i> 03/22/2018
<i>Author:</i> Stephanie Parker	<i>Revision:</i> A

## CHANGE RECORD

<b>REVISION</b>	<b>DATE</b>	<b>DESCRIPTION OF CHANGE</b>
A	3/07/2018	Initial Release

## TABLE OF CONTENTS

<b>1 DESCRIPTION</b>	<b>1</b>
1.1 Purpose . . . . .	1
1.2 Scope . . . . .	1
<b>2 RELATED DOCUMENTS</b>	<b>2</b>
2.1 Associated Documents . . . . .	2
<b>3 DATA PRODUCT DESCRIPTION</b>	<b>3</b>
3.1 Spatial Sampling Design . . . . .	3
3.2 Temporal Sampling Design . . . . .	6
3.3 Variables Reported . . . . .	7
3.4 Temporal Resolution and Extent . . . . .	7
3.5 Spatial Resolution and Extent . . . . .	7
3.6 Associated Data Streams . . . . .	8
3.7 Product Instances . . . . .	8
3.8 Data Relationships . . . . .	8
3.9 Special Considerations . . . . .	10
3.9.1 Retrieving Metabarcoding Sequence Data . . . . .	10
<b>4 DATA QUALITY</b>	<b>11</b>
4.1 Data Entry Constraint and Validation . . . . .	11
4.2 Automated Data Processing Steps . . . . .	12
4.3 Sequencing Data . . . . .	12
4.4 Data Revision . . . . .	12
4.5 Quality Flagging . . . . .	12
4.6 Analytical Facility Data Quality . . . . .	12
<b>5 REFERENCES</b>	<b>14</b>

## LIST OF TABLES AND FIGURES

Figure 1	Generic aquatic site layout for lakes, river and wadeable streams, with macroinvertebrate sampling locations in red. . . . .	5
Figure 2	Generic aquatic site layout for lakes with zooplankton sampling locations in red. . . . .	6
Figure 3	Schematic of the applications used by field technicians to enter zooplankton field data. Boxes with a gray header indicate general information that is filled out one time per sampling event (bout), and applies to all of the subsequent data. The boxes in green indicate data that are populated for each stream transect and point collected. . . . .	13

# 1 DESCRIPTION

## 1.1 Purpose

This document provides an overview of the data included in this NEON Level 1 data product, the quality controlled product generated from raw Level 0 data, and associated metadata. In the NEON data products framework, the raw data collected in the field, for example, the macroinvertebrate or zooplankton DNA samples collected in the field are considered the lowest level (Level 0). Raw data that have been quality checked via the steps detailed herein, as well as simple metrics that emerge from the raw data are considered Level 1 data products.

The text herein provides a discussion of measurement theory and implementation, data product provenance, quality assurance and control methods used, and approximations and/or assumptions made during L1 data creation.

## 1.2 Scope

This document describes the steps needed to generate the L1 data product Zooplankton collection and associated metadata from input data. This document also provides details relevant to the publication of the data products via the NEON data portal, with additional detail available in the files NEON Data Variables for Macroinvertebrate DNA Barcode (NEON.DP1.20126) (AD[05]) and NEON Data Variables for Zooplankton DNA Barcode (NEON.DP1.20221) (AD[06]) provided in the download package for this data product.

This document describes the process for ingesting and performing automated quality assurance and control procedures on the data collected in the field pertaining to AOS Protocol and Procedure: Aquatic Macroinvertebrate Sampling (AD[08]) and AOS Protocol and Procedure: Zooplankton Sampling in Lakes (AD[09]). The raw data that are processed in this document are detailed in the files NEON Raw Data Validation for Macroinvertebrate DNA Barcode (NEON.DP0.20126) (AD[03]) and NEON Raw Data Validation for Zooplankton DNA Barcode (NEON.DP0.20221) (AD[04]), provided in the download package for this data product. Please note that raw data products (denoted by ‘DP0’) may not always have the same numbers (e.g., ‘20221’) as the corresponding L1 data product.

## 2 RELATED DOCUMENTS

### 2.1 Associated Documents

AD[01]	NEON.DOC.000001	NEON Observatory Design (NOD) Requirements
AD[02]	NEON.DOC.002652	NEON Level 1, Level 2 and Level 3 Data Products Catalog
AD[03]	NEON.DP0.20126.001_dataValidation.csv	NEON Raw Data Validation for Macroinvertebrate DNA Barcode (NEON.DP0.20126)
AD[04]	NEON.DP0.20221.001_dataValidation.csv	NEON Raw Data Validation for Zooplankton DNA Barcode (NEON.DP0.20221)
AD[05]	NEON.DP1.20126.001_variables.csv	NEON Data Variables for Macroinvertebrate DNA Barcode (NEON.DP1.20126)
AD[06]	NEON.DP1.20221.001_variables.csv	NEON Data Variables for Zooplankton DNA Barcode (NEON.DP1.20221)
AD[07]	NEON.DOC.001152	NEON Aquatic Sampling Strategy
AD[08]	NEON.DOC.003046	AOS Protocol and Procedure: Aquatic Macroinvertebrate Sampling
AD[09]	NEON.DOC.001194	AOS Protocol and Procedure: Zooplankton Sampling in Lakes
AD[10]	NEON.DOC.000008	NEON Acronym List
AD[11]	NEON.DOC.000243	NEON Glossary of Terms
AD[12]	OS_Generic_Transitions.pdf	NEON Algorithm Theoretical Basis Document: OS Generic Transitions
AD[13]	Nicl Language.pdf	NEON's Ingest Conversion Language (NICL) specifications

### 3 DATA PRODUCT DESCRIPTION

The DNA metabarcoding data products provide DNA sequence data and metadata for macroinvertebrate and zooplankton communities at NEON aquatic sites. These data may be used for determination of diversity patterns in bulk samples, and analogous morphological taxonomy samples for both macroinvertebrates and zooplankton are collected at the same time and location, so taxonomic data may be correlated by data users. NEON uses PCR amplification of two target regions of the CO1 gene using primers described in Gibson et al. (2015). The use of two distinct primer sets and regions enables greater coverage of the diversity of distinct arthropod groups. Sequence data are generated using high-throughput technology that produces many thousands of sequence reads per sample (Armougom and Didier 2009, Klindworth et al. 2013).

The type of sampler used to collect a sample in the field is determined by the habitat and substrate type (macroinvertebrates) or water depth (zooplankton) at the sampling location. Macroinvertebrates are collected using a Surber sampler, modified kicknet, hand corer, or D-frame net (lakes and rivers only). All sampling devices collect material from a known area of the benthos. For zooplankton samples, locations deeper than 4 m are sampled using a vertical tow net, while locations shallower than 4 m are sampled using a Schindler-Patalas sampler (USEPA 2012a, 2012b). Typically, multiple (up to 3) tows or Schindler traps are collected and composited into a single sample. Zooplankton samples are collected on a volumetric basis.

Sample collection methods differ between macroinvertebrate and zooplankton samples, but in general samples are minimally processed in the field in order to reduce the introduction of microbial contaminants. After collection, samples are preserved in 95% ethanol and chilled at 4 degrees C, with an ethanol change occurring within 48 hours of collection. For additional information see sampling design NEON Aquatic Sampling Strategy (AD[07]), and protocols AOS Protocol and Procedure: Aquatic Macroinvertebrate Sampling (AD[08]) and AOS Protocol and Procedure: Zooplankton Sampling in Lakes (AD[09]). Samples are shipped to an analytical laboratory where sample homogenization, DNA extraction, sequencing library preparation, and DNA sequencing occur.

#### 3.1 Spatial Sampling Design

Benthic invertebrates at NEON aquatic sites (Figure 1) are sampled using a percent-based macrohabitat approach (after Moulton et al. 2002). Habitats sampled focus on riffles, runs, pools, and step pools depending on the percent cover of each habitat within each 1 km-long NEON Aquatic wadeable stream site (NOTE: some NEON sites may be less than 1 km due to permitting restrictions), and benthic-pelagic and littoral samples in lakes and non-wadeable streams. Three samples are collected in the dominant habitat type (wadeable stream) or littoral area (lake and non-wadeable stream) on a given sampling date at a site.

Samplers used for macroinvertebrate collection are designed to work by disturbing the benthic sediments and catching invertebrates in an attached net or container, while delineating the benthic area sampled for a quantitative result. The sampler type chosen differs depending on the water depth, velocity, and substratum type in the chosen habitat (Hauer and Resh 2006). The

<i>Title:</i> NEON User Guide to Macroinvertebrate DNA Barcode (NEON.DP1.20126) and Zooplankton DNA Barcode (NEON.DP1.20221)	<i>Date:</i> 03/22/2018
<i>Author:</i> Stephanie Parker	<i>Revision:</i> A

collection method may differ depending on the habitat and substrate being sampled, however all samples are collected from the surface of the natural substratum in each habitat using a quantitative sampling method. See AOS Protocol and Procedure: Aquatic Macroinvertebrate Sampling (AD[08]) for additional details on sampling strategy and SOPs.

Zooplankton samples are collected only from lakes (Figure 2). Samples are collected near the inlet, outlet, and buoy (deepest) sampling locations and are designed to sample organisms inhabiting the water column. The type of sampler selected depends on the depth of the lake at the sampling location, and the volume of lake water sampled can be calculated using the number of tows/traps and the volume of the sampler used. See AOS Protocol and Procedure: Zooplankton Sampling in Lakes (AD[09]) for additional details on sampling strategy and SOPs.

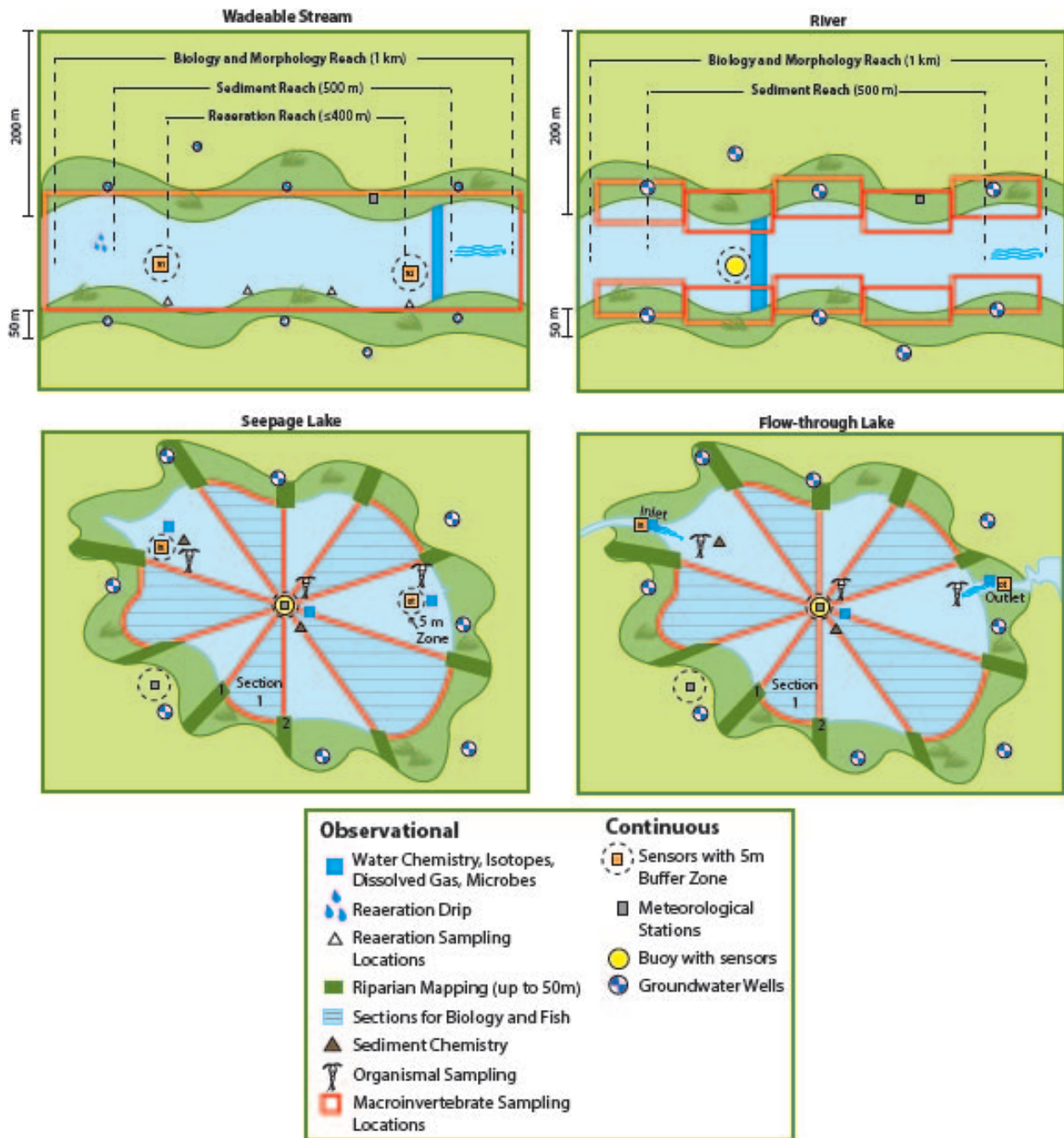


Figure 1: Generic aquatic site layout for lakes, river and wadeable streams, with macroinvertebrate sampling locations in red.



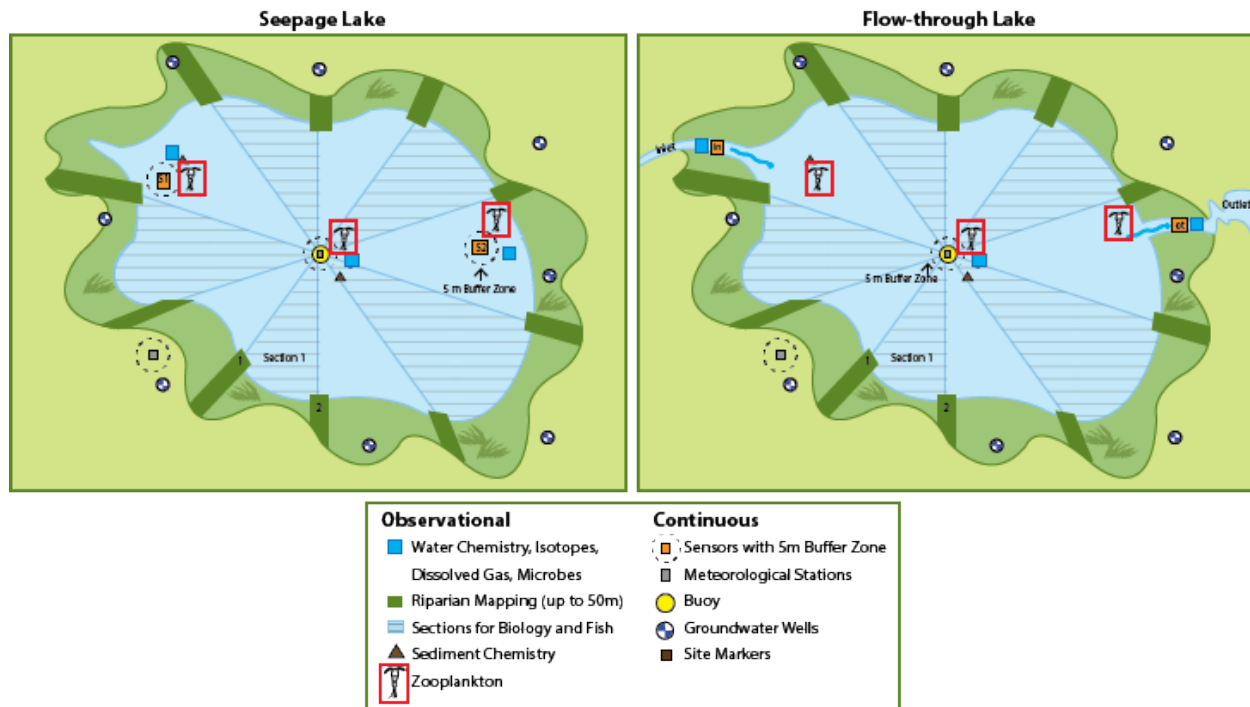


Figure 2: Generic aquatic site layout for lakes with zooplankton sampling locations in red.

### 3.2 Temporal Sampling Design

Sampling occurs three times per year at each NEON site (AD[07]). Timing of sampling is site-specific and determined based on historical hydrological and meteorological data. Sample bout 1 is an early-season date, representing a period of rapid biomass accumulation after winter, typically prior to leaf out or after ice-off where applicable. Sample bout 2 targets mid-summer base-flow conditions and sample bout 3 represents the late growing season (typically autumn) during leaf-fall where applicable. These dates differ on a site-by-site basis, but should always occur at, or near, baseflow conditions within the watershed. Sampling does not occur directly following a rain or wind event that causes turbidity in the water column (lakes/rivers) or a flood in wadeable streams (defined as  $>1.5 \times$  base flow; Biggs et al. 1999). Should such a flood event occur on or prior to a target collection date, sampling is delayed 3 days-1 week (maximum 2 weeks, dependent on field schedule) to allow for invertebrates to recolonize the substratum (c.f. Brooks and Boulton 1991, Matthaei et al. 1996). Sampling at each site is completed within a single day for each bout. See NEON Aquatic Sampling Strategy (AD[07]), AOS Protocol and Procedure: Aquatic Macroinvertebrate Sampling (AD[08]), and AOS Protocol and Procedure: Zooplankton Sampling in Lakes (AD[09]) for additional details.

### 3.3 Variables Reported

All variables reported from the field or laboratory technician (L0 data) are listed in the files NEON Raw Data Validation for Macroinvertebrate DNA Barcode (NEON.DP0.20126) (AD[03]) and NEON Raw Data Validation for Zooplankton DNA Barcode (NEON.DP0.20221) (AD[04]). All variables reported in the published data (L1 data) are also provided separately in the files NEON Data Variables for Macroinvertebrate DNA Barcode (NEON.DP1.20126) (AD[05]) and NEON Data Variables for Zooplankton DNA Barcode (NEON.DP1.20221) (AD[06]).

Field names have been standardized with Darwin Core terms (<http://rs.tdwg.org/dwc/>; accessed 7 December 2017), the Global Biodiversity Information Facility vocabularies (<http://rs.gbif.org/vocabulary/gbif/>; accessed 7 December 2017), the VegCore data dictionary (<https://projects.nceas.ucsb.edu/nceas/projects/bien/wiki/VegCore>; accessed 7 December 2017), where applicable. NEON Aquatic Observation System (AOS) spatial data employs the World Geodetic System 1984 (WGS84) for its fundamental reference datum and Earth Gravitational Model 96 (EGM96) for its reference gravitational ellipsoid. Latitudes and longitudes are denoted in decimal notation to six decimal places, with longitudes indicated as negative west of the Greenwich meridian.

Some variables described in this document may be for NEON internal use only and will not appear in downloaded data.

### 3.4 Temporal Resolution and Extent

The finest temporal resolution that macroinvertebrate and zooplankton DNA data will be tracked is per sampling day. All 3 samples per module (macroinvertebrate or zooplankton) are collected within a single day at a particular site. A suite of other biological sampling occurs at the site during the same ~30 day bout. Three sampling bouts occur per site per year.

The finest resolution at which temporal data are reported is at **collectDate**, the date and time of day when the samples were collected in the field.

The NEON Data Portal provides data in monthly files for query and download efficiency. Queries including any part of a month will return data from the entire month. Code to stack files across months is available here: <https://github.com/NEONScience/NEON-utilities>.

### 3.5 Spatial Resolution and Extent

Each macroinvertebrate sample represents a patch of stream bottom within the 1 km permitted wadeable or non-wadeable stream reach, or permitted lake area, and contains multiple individuals. The exact location (latitude and longitude) of each sample is not tracked as it is intended to represent the overall habitat. The **namedLocation** reported in wadeable streams represents a midpoint in the permitted reach, plus coordinate uncertainty surrounding that point. In lakes and non-wadeable streams, some samples are collected near monumented locations associated with a more-specific **namedLocation** (e.g., NEON sensor infrastructure). Sampling locations are tracked by latitude and longitude and include an indication of **coordinateUncertainty**.

Each zooplankton sample represents a location in a lake near one of the NEON sensor installations (inlet, outlet, or buoy), and contains multiple individuals. The **namedLocation** reported represents the location of the NEON sensor infrastructure near the sampling location, plus coordinate uncertainty surrounding that point. The protocol dictates that samples are collected approximately 5 m from the sensor infrastructure to minimize effects on the profiling data, so the standard **coordinateUncertainty** is 10 m to represent the normal sampling distance from the sampling location. If, for some reason, sampling cannot occur within 10 m of the named location, technicians will enter **additionalCoordinateUncertainty**.

Samples are collected from the dominant habitat type (wadeable streams), benthic littoral zone (lakes/rivers), or pelagic water column (lakes - zooplankton). Overall, this results in a spatial hierarchy of:

namedLocation (finest spatial resolution, ID of location within site) -> siteID (ID of NEON site)  
-> domainID (ID of a NEON domain)

### 3.6 Associated Data Streams

Macroinvertebrate and zooplankton DNA metabarcode samples are collected at the same time and location, and using the same method, as an analogous morphological taxonomy sample. Related samples share the same **eventID** and **namedLocation**, as well as the same root **sampleID**. DNA sample IDs are equal to the taxonomic sampleID + “DNA” appended to the end. Taxonomic data are available in the NEON data products “Macroinvertebrate Collection” (DP1.20120.001) and “Zooplankton Collection” (DP1.20219.001).

### 3.7 Product Instances

At each aquatic site, there will be up to 9 samples macroinvertebrate samples collected per year (3 macroinvertebrate samples per sampling bout). At lake sites, there will be up to 9 zooplankton samples collected per year. Each sample may generate multiple records from the external lab.

### 3.8 Data Relationships

#### 3.8.0.1 Macroinvertebrate DNA Barcode (NEON.DP1.20120)

For each macroinvertebrate sampling event, a record is created in `inv_fieldData`. In the event that sampling is impractical (e.g., the location is dry, ice covered, etc.) or no **geneticSampleID** is taken, and there will be no child records. Otherwise, there may be a number of child records in subsequent tables. Child records will be found in `inv_dnaExtraction` (initial subsampling and dna extraction metadata at the external facility), `inv_pcrAmplification` (PCR metadata), and `inv_markerGeneSequencing` (sequencing metadata). Every record in `inv_dnaExtraction`, `inv_pcrAmplification`, and `inv_markerGeneSequencing` should have a corresponding record in `inv_fieldData` describing field collection conditions, location, and metadata during sample collection. The **dnaSampleID** is created in the `inv_dnaExtraction` table, linking downstream data

<i>Title:</i> NEON User Guide to Macroinvertebrate DNA Barcode (NEON.DP1.20126) and Zooplankton DNA Barcode (NEON.DP1.20221)	<i>Date:</i> 03/22/2018
<i>Author:</i> Stephanie Parker	<i>Revision:</i> A

to the **geneticSampleID** from the `inv_fieldData` table. There is one unique record for each **dnaSampleID** in `inv_dnaExtraction` unless extractions were unsuccessful and multiple extractions were required, while `inv_pcrAmplification` and `inv_markerGeneSequencing` may have multiple records per **dnaSampleID**, corresponding to different **replicates**. Duplicates and/or missing data may exist where protocol and/or data entry aberrations have occurred; users should check data carefully for anomalies before joining tables.

`inv_fieldData.csv` -> One record is created for each sample collected in the field, creating a unique **geneticSampleID** representing one sample per **collectDate** and **namedLocation**. This table also indicates field conditions, including **samplerType**, **habitatType**, and **benthicArea**.

`inv_dnaExtraction.csv` -> The **geneticSampleID** from the `fieldData` table is subsampled to create the **dnaSampleID** in the `inv_dnaExtraction` table. Generally, there will be only one DNA extraction per **dnaSampleID**, but in some cases multiple extractions will be necessary.

`inv_pcrAmplification.csv` -> Metadata on PCR amplification sample processing is presented in this table. One record is expected per **dnaSampleID/replicate** combination. The **dnaSampleID** equals the **dnaSampleID** in the `inv_dnaExtraction` table.

`inv_markerGeneSequencing.csv` -> Metadata on gene sequencing is presented in this table. One record is expected per **dnaSampleID/replicate** combination. The **dnaSampleID** equals the **dnaSampleID** in the `inv_dnaExtraction` table.

### 3.8.0.2 Zooplankton DNA Barcode (NEON.DP1.20221)

For each zooplankton sampling event where a genetic sample is collected, a record is created in `zoo_fieldData` with a unique **geneticSampleID**. In the event that sampling is impractical (e.g., the location is dry, ice covered, etc.) or no **geneticSampleID** is taken, no `zoo_fieldData` records will appear in the download. If sampling occurs, there may be a number of child records in subsequent tables. Child records will be found in `zoo_dnaExtraction` (initial subsampling and dna extraction metadata at the external facility), `zoo_pcrAmplification` (PCR metadata), and `zoo_markerGeneSequencing` (sequencing metadata). Every record in `zoo_dnaExtraction`, `zoo_pcrAmplification`, and `zoo_markerGeneSequencing` should have a corresponding record in `zoo_fieldData` describing field collection conditions, location, and metadata during sample collection. The **dnaSampleID** is created in the `zoo_dnaExtraction` table, linking to **geneticSampleID** from the `zoo_fieldData` table (`zoo_dnaExtraction.geneticSampleID=zoo_fieldData.sampleID`). There is one unique record for each **dnaSampleID** in `zoo_dnaExtraction` unless extractions were unsuccessful and multiple extractions were required, while `zoo_pcrAmplification` and `zoo_markerGeneSequencing` may have multiple records per **dnaSampleID** corresponding to different **replicates**. Duplicates and/or missing data may exist where protocol and/or data entry aberrations have occurred; users should check data carefully for anomalies before joining tables.

`zoo_fieldData.csv` -> One record is created for each sample collected in the field, creating a unique **geneticSampleID**. This table also indicates the field conditions, including **sam-**

<i>Title:</i> NEON User Guide to Macroinvertebrate DNA Barcode (NEON.DP1.20126) and Zooplankton DNA Barcode (NEON.DP1.20221)	<i>Date:</i> 03/22/2018
<i>Author:</i> Stephanie Parker	<i>Revision:</i> A

**plerType**, number of tows or traps collected (**towsTrapsNumber**), and sampling depth (**zooDepth1**, **zooDepth2**, **zooDepth3**).

zoo\_dnaExtraction.csv - > The **geneticSampleID** from the fieldData table is subsampled to create the **dnaSampleID** in the inv\_dnaExtraction table. One record is expected per **dnaSampleID** here, and will be linked to subsequent tables. Generally, there will be only one DNA extraction per **dnaSampleID**, but in some cases multiple extractions will be necessary.

zoo\_pcrAmplification.csv - > Metadata on PCR amplification sample processing is presented in this table. One record is expected per **dnaSampleID/replicate** combination. The **dnaSampleID** equals the **dnaSampleID** in the inv\_dnaExtraction table.

zoo\_markerGeneSequencing.csv - > Metadata on gene sequencing is presented in this table One record is expected per **dnaSampleID/replicate** combination. The **dnaSampleID** equals the **dnaSampleID** in the inv\_dnaExtraction table.

### 3.9 Special Considerations

For ease of integration with external data sets, metabarcode sequence data are published on public sequence repositories. The primary data repository is MG-RAST (<http://metagenomics.anl.gov>, Meyer et al., 2008), which directly synchronizes its data with the European Bioinformatics Institute (EMBL-EBI) database and, through EMBL-EBI, synchronizes with the National Center for Biotechnology Information’s Sequence Read Archive (SRA). A suite of metadata, compliant with minimum metadata standards defined by the Genomics Standards Consortium (e.g. MIxS, MIMARKS), accompanies the sequence data. While efforts are made to publish comprehensive sequencing metadata with the sequence data stored at public sequence repositories, potentially important data will only be available through the NEON Data Portal. These data include:

- Methods and SOPs
- QA data
- Sample identifiers to enable joining metabarcoding data with other related Data Products
- Data for other related Data Products

The sequence data for each region of the CO1 gene amplified are uploaded separately, such that there will be one file per sample per gene region targeted.

#### 3.9.1 Retrieving Metabarcoding Sequence Data

There are a number of ways to search and retrieve minimally processed metabarcoding sequence data.

- From the NEON data portal:
  1. Links beginning with “MG-RAST Project: NEON Macroinvertebrate DNA Barcode” will take the user to the MG-RAST project page for the queried data. This is a dynamic link and will automatically update based on the user query.

<i>Title:</i> NEON User Guide to Macroinvertebrate DNA Barcode (NEON.DP1.20126) and Zooplankton DNA Barcode (NEON.DP1.20221)	<i>Date:</i> 03/22/2018
<i>Author:</i> Stephanie Parker	<i>Revision:</i> A

2. The link “MG-RAST Project: NEON Zooplankton DNA Barcode” will take the user to the MG-RAST project page for the queried data. This is a dynamic link and will automatically update based on the user query.
  3. The link “MG-RAST Sample Search” takes the user to the MG-RAST page for searching individual records, pre-populated with NEON records based on the user query.
- From MG-RAST directly: Users who are interested in using the MG-RAST data analysis pipeline may want to combine NEON datasets with other datasets. This may be more easily achieved by querying the MG-RAST database directly. Users can analyze samples from a variety of NEON and non-NEON projects. A free user account may be required.
  - From SRA directly: Data and metadata are available for download from the SRA using the SRA toolkit. Documentation on how to install and use the toolkit for downloading sequence data is available on the SRA website.
  - From EMBL-EBI: MG-RAST also synchronizes data sets with the European Bioinformatics Initiative Repository (EMBL-EBI, <https://www.ebi.ac.uk/>), which has a web and API interface for downloading data. The NEON macroinvertebrate metabarcode data can be found by querying the NCBI Project ID PRJNA391345, and the NEON zooplankton metabarcode data can be found by querying the NCBI Project ID PRJNA391744.

*Note:* There may be lags between publication of metadata on the NEON data portal and availability of sequence data on the public sequence repository.

## 4 DATA QUALITY

### 4.1 Data Entry Constraint and Validation

Many quality control measures are implemented at the point of data entry within a mobile data entry application or web user interface (UI). For example, data formats are constrained and data values controlled through the provision of dropdown options, which reduces the number of processing steps necessary to prepare the raw data for publication. The field data entry workflow for collecting zooplankton field data is diagrammed in Figure 3.

An additional set of constraints are implemented during the process of ingest into the NEON database. The product-specific data constraint and validation requirements built into data entry applications and database ingest are described in the documents NEON Raw Data Validation for Macroinvertebrate DNA Barcode (NEON.DP0.20126) (AD[03]) and NEON Raw Data Validation for Zooplankton DNA Barcode (NEON.DP0.20221) (AD[04]), provided with every download of this data product. Contained within this file is a field named ‘entryValidationRulesForm’, which describes syntactically the validation rules for each field built into the data entry application. Data entry constraints are described in NiCl syntax in the validation file provided with every data download, and the NiCl language is described in NEON’s Ingest Conversion Language (NICL) specifications ([AD[13]]).

Data collected prior to 2017 were processed using a paper-based workflow that did not implement the full suite of quality control features associated with the interactive digital workflow.

## 4.2 Automated Data Processing Steps

### 4.3 Sequencing Data

Sequencing data are generated in batches of multiple samples. After sequencing, the multiplexed sequence data are parsed into separate files on a per sample basis. For each sample, minimum quality criteria must be met in order to accept the data for the sample. The general criteria include meeting a minimum sequencing depth (e.g. number of sequences per sample), a maximum number of ambiguous base calls, and a minimum quality score. The actual criteria may change over time as technology evolves and standards change. The per sample QA results are published as part of the expanded download package.

Following data entry into a mobile application of web user interface, the steps used to process the data through to publication on the NEON Data Portal are detailed in the NEON Algorithm Theoretical Basis Document: OS Generic Transitions (AD[12]).

### 4.4 Data Revision

All data are provisional until a numbered version is released; the first release of a static version of NEON data, annotated with a globally unique identifier, is planned to take place in 2020. During the provisional period, QA/QC is an active process, as opposed to a discrete activity performed once, and records are updated on a rolling basis as a result of scheduled tests or feedback from data users. The Change Log section of the data product readme, provided with every data download, contains a history of major known errors and revisions.

### 4.5 Quality Flagging

The **dataQF** field in each data record is a quality flag for known errors applying to the record. Please see the table below for an explanation of **dataQF** codes specific to this product.

fieldName	value	definition
dataQF	legacyData	Data recorded using a paper-based workflow that did not implement the full suite of quality control features associated with the interactive digital workflow

### 4.6 Analytical Facility Data Quality

Data analyses conducted on sequencing data conform to the current data quality standards used by practitioners. Each metadata table includes a variable, called **qaqcStatus**, in which the labo-

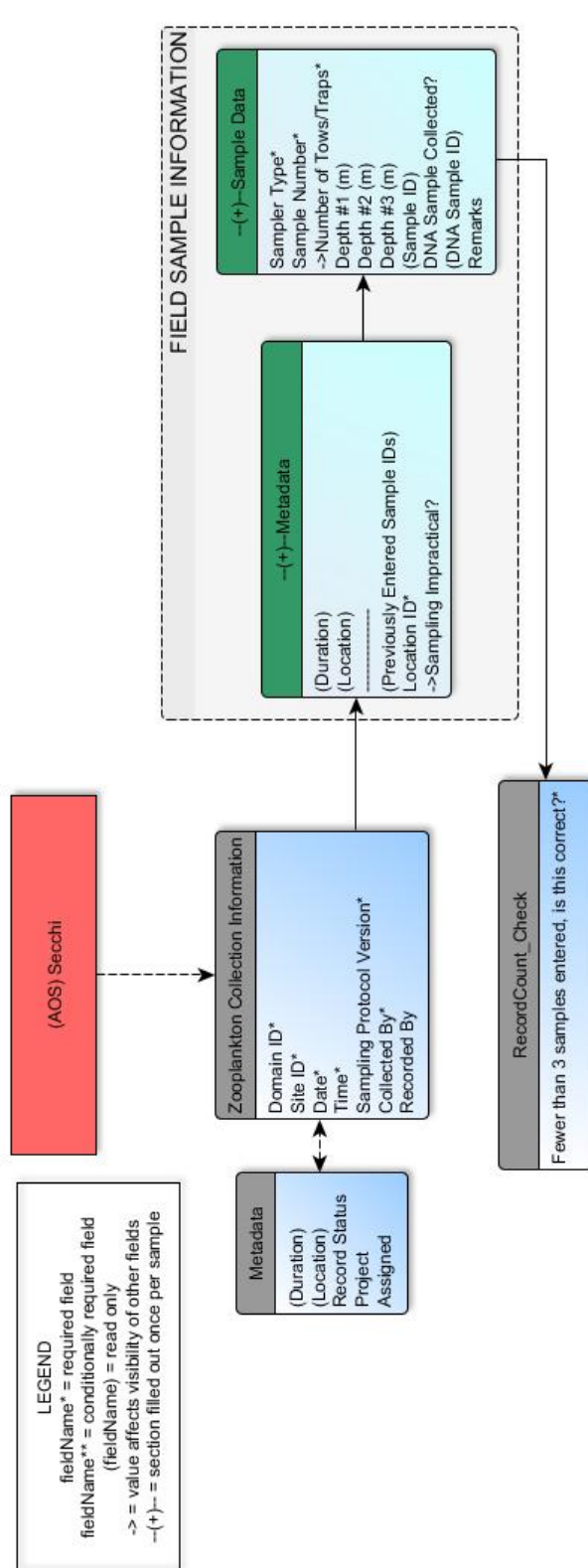


Figure 3: Schematic of the applications used by field technicians to enter zooplankton field data. Boxes with a gray header indicate general information that is filled out one time per sampling event (bout), and applies to all of the subsequent data. The boxes in green indicate data that are populated for each stream transect and point collected.



ratory can indicate sample processing issues. Any records with a qaqcStatus = “Fail” should also be accompanied by free-form notes in the “remarks” variable.

## 5 REFERENCES

- Armougom F., and R. Didier. 2009. Exploring microbial diversity using 16S rRNA high-throughput methods. *Journal of Computer Science and Systems Biology* 2:74-92. <https://doi.org/10.4172/jcsb.1000019>.
- Biggs, B. J. F., R. A. Smith, and M. J. Duncan. 1999. Velocity and sediment disturbance of periphyton in headwater streams: biomass and metabolism. *Journal of the North American Benthological Society* 18: 222-241.
- Brooks, S. S. and A. J. Boulton. 1991. Recolonization dynamics of benthic macroinvertebrates after artificial and natural disturbances in an Australian temporary stream. *Australian Journal of Marine and Freshwater Research* 42:295-308.
- Gibson J. F., S. Shokralla, C. Curry, D. J. Baird, W. A. Monk, I. King, and M. Hajibabaei. 2015. Large-scale biomonitoring of remote and threatened ecosystems via high-throughput sequencing. *PLoS One*. 10: e0138432.
- Hauer, F. R. and V. H. Resh. 2006. Macroinvertebrates. Pages 435-463 in F. R. Hauer and G. A. Lamberti, editors. *Methods in Stream Ecology*, Second Edition. Academic Press, Boston, MA.
- Klindworth A., E. Pruesse, T. Schweer, J. Peplies, C. Quast, M. Horn, and F. O. Glockner. 2013. Evaluation of general 16S ribosomal RNA gene PCR primers for classical and next-generation sequencing-based diversity studies. *Nucleic Acids Research* 41:e1-e1.
- Matthaei, C. D., U. Uhlinger, E. I. Meyer, and A. Frutiger. 1996. Recolonization by benthic invertebrates after experimental disturbance in a Swiss prealpine river. *Freshwater Biology* 35: 233-248.
- Meyer F., D. Paarmann, M. D’Souza, R. Olson, E.M. Glass, M. Kubal, T. Paczian, et al. 2008. The Metagenomics RAST Server - a public resource for the automatic phylogenetic and functional analysis of metagenomes. *BMC Bioinformatics* 9: 386.
- Moulton, S. R., II, J. G. Kennen, R. M. Goldstein, and J. A. Hambrook. 2002. Revised protocols for sampling algal, invertebrate, and fish communities as part of the National Water-Quality Assessment Program. Open-File Report 02-150. U.S. Geological Survey, Reston, VA.
- USEPA. 2012a. National Lakes Assessment Program, Field Operations Manual.
- USEPA. 2012b. Sampling Procedures for the Great Lakes.