

<i>Title:</i> NEON User Guide to Microbe Group Abundances (NEON.DP1.10109; NEON.DP1.20277; NEON.DP1.20278)	<i>Date:</i> 02/27/2018
<i>Author:</i> Lee Stanish	<i>Revision:</i> A

## NEON USER GUIDE TO MICROBE GROUP ABUNDANCES (NEON.DP1.10109; NEON.DP1.20277; NEON.DP1.20278)

PREPARED BY	ORGANIZATION	DATE
Lee Stanish	FSU	02/27/2018
Stephanie Parker	AOS	02/27/2018

<i>Title:</i> NEON User Guide to Microbe Group Abundances (NEON.DP1.10109; NEON.DP1.20277; NEON.DP1.20278)	<i>Date:</i> 02/27/2018
<i>Author:</i> Lee Stanish	<i>Revision:</i> A

## CHANGE RECORD

REVISION	DATE	DESCRIPTION OF CHANGE
A	11/8/2017	Initial Release

## TABLE OF CONTENTS

<b>1</b>	<b>DESCRIPTION</b>	<b>1</b>
1.1	Purpose . . . . .	1
1.2	Scope . . . . .	1
<b>2</b>	<b>RELATED DOCUMENTS AND ACRONYMS</b>	<b>2</b>
2.1	Associated Documents . . . . .	2
2.2	Acronyms . . . . .	2
<b>3</b>	<b>DATA PRODUCT DESCRIPTION</b>	<b>3</b>
3.1	Spatial Sampling Design . . . . .	4
3.2	Temporal Sampling Design . . . . .	6
3.3	Variables Reported . . . . .	7
3.4	Spatial Resolution and Extent . . . . .	7
3.4.1	Soils . . . . .	7
3.4.2	Aquatics . . . . .	8
3.5	Temporal Resolution and Extent . . . . .	8
3.6	Associated Data Streams . . . . .	8
3.6.1	Soils . . . . .	8
3.6.2	Aquatics . . . . .	9
3.7	Product Instances . . . . .	9
3.8	Data Relationships . . . . .	10
3.8.1	Soils . . . . .	10
3.8.2	Aquatics . . . . .	11
<b>4</b>	<b>DATA QUALITY</b>	<b>13</b>
4.1	Data Entry Constraint and Validation . . . . .	13
4.2	Automated Data Processing Steps . . . . .	13
4.3	Data Revision . . . . .	13
4.4	Quality Flagging . . . . .	14
4.5	Analytical Facility Data Quality . . . . .	14
<b>5</b>	<b>REFERENCES</b>	<b>14</b>

## LIST OF TABLES AND FIGURES

Figure 1	Overview of microbial field sample types, processing steps, and analyses. . . . .	4
Figure 2	Overview of soil microbial field sampling and analysis workflow. . . . .	5
Figure 3	Generic NEON aquatic site layouts with microbial sampling locations highlighted in red boxes. . . . .	6

<i>Title:</i> NEON User Guide to Microbe Group Abundances (NEON.DP1.10109; NEON.DP1.20277; NEON.DP1.20278)	<i>Date:</i> 02/27/2018
<i>Author:</i> Lee Stanish	<i>Revision:</i> A

## 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, soil temperature from a single collection event - 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 products for Microbe Group Abundances data and associated metadata measured on aquatic and terrestrial samples by broad taxonomic (e.g. bacterial, archaeal, and fungal) groups. 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 Soil Microbe Group Abundances (NEON.DP1.10109) (AD[05]), NEON Data Variables for Benthic Microbe Group Abundances (NEON.DP1.20277) (AD[06]), and NEON Data Variables for Surface Water Microbe Group Abundances (NEON.DP1.20278) (AD[07]), provided in the download package for each of these three data products.

This document describes the process for ingesting and performing automated quality assurance and control procedures on the laboratory data from samples generated by the following field sampling protocols: TOS Protocol and Procedure: Soil Biogeochemical and Microbial Sampling (AD[10]) for upland soil samples; with TOS Standard Operating Procedure: Wetland Soil Sampling (AD[11]) for wetland soil samples; or AOS Protocol and Procedure: Aquatic Microbial Sampling (AD[12]) for aquatic samples. The raw data that are processed as described in this document are detailed in the file, NEON Raw Data Validation for Microbe Group Abundances (NEON.DP0.10109) (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., '10033') as the corresponding L1 data product.

## 2 RELATED DOCUMENTS AND ACRONYMS

### 2.1 Associated Documents

AD[01]	NEON.DOC.000001	NEON Observatory Design (NOD) Requirements
AD[02]	NEON.DOC.000913	TOS Science Design for Spatial Sampling
AD[03]	NEON.DOC.002652	NEON Level 1, Level 2 and Level 3 Data Products Catalog
AD[04]	NEON.DP0.10109.001_dataValidation.csv	NEON Raw Data Validation for Microbe Group Abundances (NEON.DP0.10109)
AD[05]	NEON.DP1.10109.001_variables.csv	NEON Data Variables for Soil Microbe Group Abundances (NEON.DP1.10109)
AD[06]	NEON.DP1.20277.001_variables.csv	NEON Data Variables for Benthic Microbe Group Abundances (NEON.DP1.20277)
AD[07]	NEON.DP1.20278.001_variables.csv	NEON Data Variables for Surface Water Microbe Group Abundances (NEON.DP1.20278)
AD[08]	NEON.DOC.000908	TOS Science Design for Microbial Diversity
AD[09]	NEON.DOC.001152	NEON Aquatic Sample Strategy Document
AD[10]	NEON.DOC.014048	TOS Protocol and Procedure: Soil Biogeochemical and Microbial Sampling
AD[11]	NEON.DOC.003044	AOS Protocol and Procedure: Aquatic Microbial Sampling
AD[12]	NEON.DOC.000913	TOS Science Design for Spatial Sampling
AD[13]	NEON.DOC.000008	NEON Acronym List
AD[14]	NEON.DOC.000243	NEON Glossary of Terms
AD[15]	OS_Generic_Transitions.pdf	NEON Algorithm Theoretical Basis Document: OS Generic Transitions
AD[16]		NEON's Ingest Conversion Language (NICL) specifications

### 2.2 Acronyms

Acronym	Definition
qPCR	Quantitative Polymerase Chain Reaction

<p><i>Title:</i> NEON User Guide to Microbe Group Abundances (NEON.DP1.10109; NEON.DP1.20277; NEON.DP1.20278)</p>	<p><i>Date:</i> 02/27/2018</p>
<p><i>Author:</i> Lee Stanish</p>	<p><i>Revision:</i> A</p>

### 3 DATA PRODUCT DESCRIPTION

The Microbe Group Abundances data products provide semi-quantitative estimates of the abundances of bacteria, archaea and fungi in soil and aquatic samples. Group abundances are quantified using Quantitative Polymerase Chain Reaction (qPCR), a method of measuring the abundance of a target gene in a sample, which is used to approximate the abundance of the organisms containing that gene within a sample. NEON measures the abundance of the 16S rRNA gene to quantify the abundances of bacteria and archaea and measures the ITS rRNA gene to quantify the abundances of fungi (Ginzinger 2002). The sample plan implements the guidelines and requirements in the Science Designs for TOS Terrestrial Microbial Diversity (AD[08]) and Aquatic Sampling (AD[09]). Information on sample collection methods such as frequencies per sample type can be found in the field user guides for each data product:

- Soils: NEON User Guide to Soil Physical Properties, Distributed Periodic (NEON.DP1.10086)
- Surface water: NEON User Guide for Surface Water Microbe Cell Count (NEON.DP1.20138)
- Benthic habitats: NEON User Guide for Aquatic Benthic Microbe Collection (NEON.DP0.20270)

In general, samples are minimally processed in the field in order to reduce the introduction of microbial contaminants. After collection, samples are frozen in the field on dry ice and transported to ultra-low freezers at the NEON field laboratories. Samples are shipped to an analytical laboratory where sample processing and qPCR analysis occurs (Figure 1). For data generated prior to Jan 1, 2016, 3 separate primer sets were used to quantify bacterial, archaeal, and fungal abundances. For data generated after Jan 1, two primer sets were used: 1 primer set that amplifies both bacteria and archaea, and 1 set that amplifies fungi. For specific methods and primer sets used, refer to the *mga\_labSummary* data table, included in this download package.

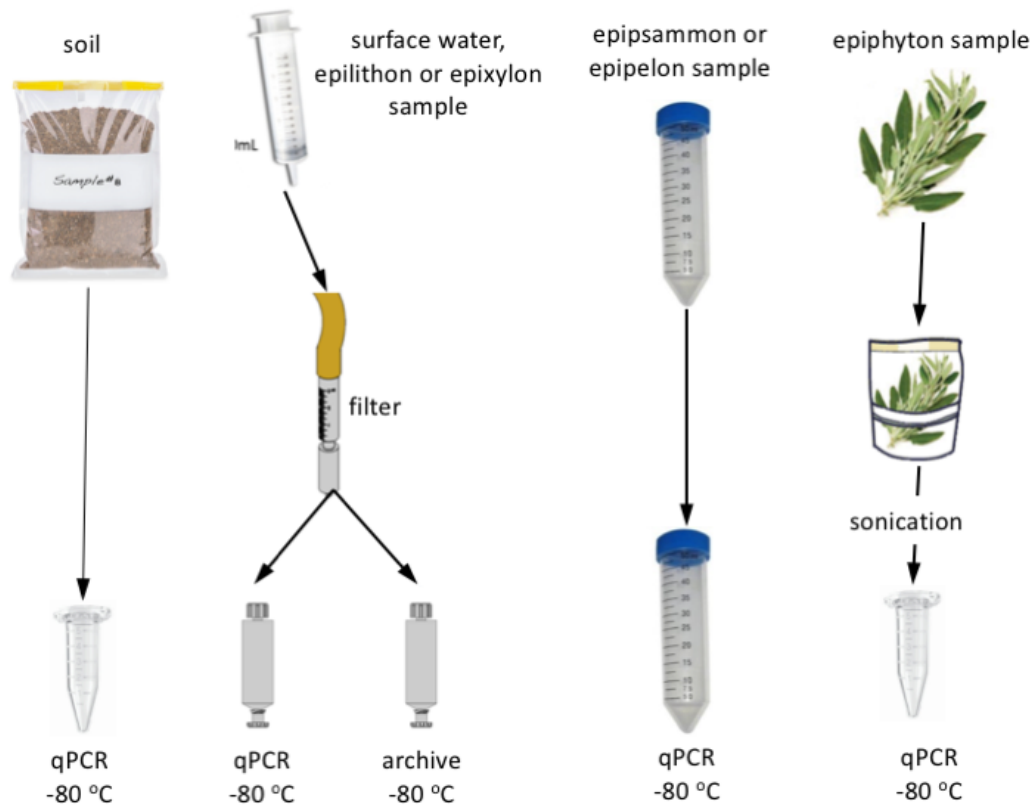


Figure 1: Overview of microbial field sample types, processing steps, and analyses.

### 3.1 Spatial Sampling Design

Sampling for microbial group abundance analysis is executed at all NEON sites and for all samples, data are reported at the resolution of a single sampling location.

For soils, this equates to a randomly-assigned X,Y coordinate (+/- 0.5 meters) within a NEON plot. Ten plots are sampled at 3 randomly selected locations within each plot (Figure 2). In general, only the surface horizon is sampled to a maximum depth of 30cm, and horizons are broadly defined as either organic (O) or mineral (M).

<p>Title: NEON User Guide to Microbe Group Abundances (NEON.DP1.10109; NEON.DP1.20277; NEON.DP1.20278)</p>	<p>Date: 02/27/2018</p>
<p>Author: Lee Stanish</p>	<p>Revision: A</p>

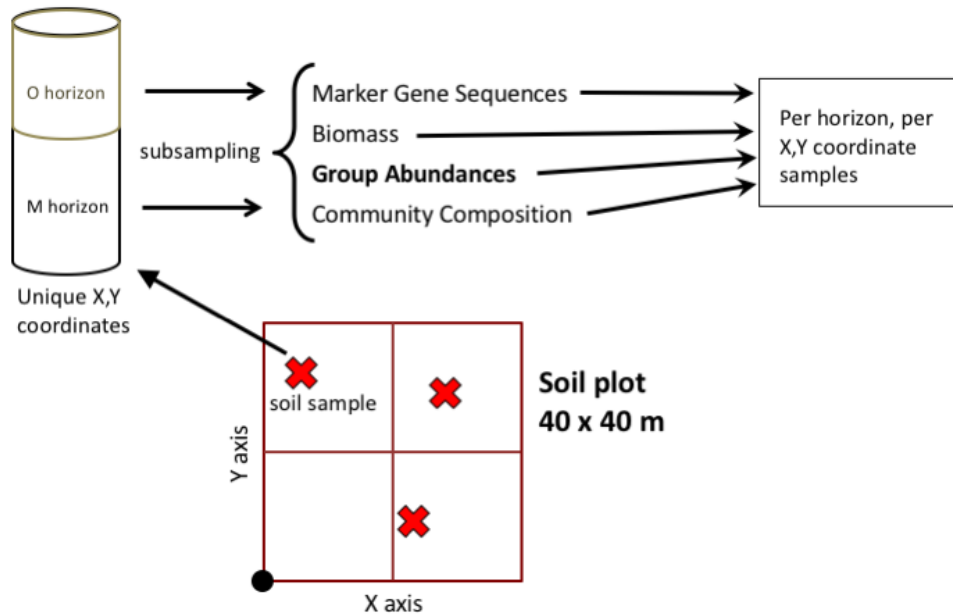


Figure 2: Overview of soil microbial field sampling and analysis workflow.

For aquatic surface water samples, this equates to the buoy sensor station and inlet/outlet locations within a lake, the buoy sensor station for large rivers, or the downstream sensor array for wadeable streams. For aquatic benthic samples, this equates to up to eight locations within a 1 km reach (Figure 3).

The spatial designs for the microbe group abundances data products are described in more detail in the Data Product User Guides for Soil Physical Properties (NEON.DP1.10086), Aquatic Surface Water Cell Counts (NEON.DP1.20138), and Aquatic Benthic Field Sampling (NEON.DP0.20270). For a description of the methods used in terrestrial plot selection, refer to the TOS Science Design for Spatial Sampling (AD[02]).



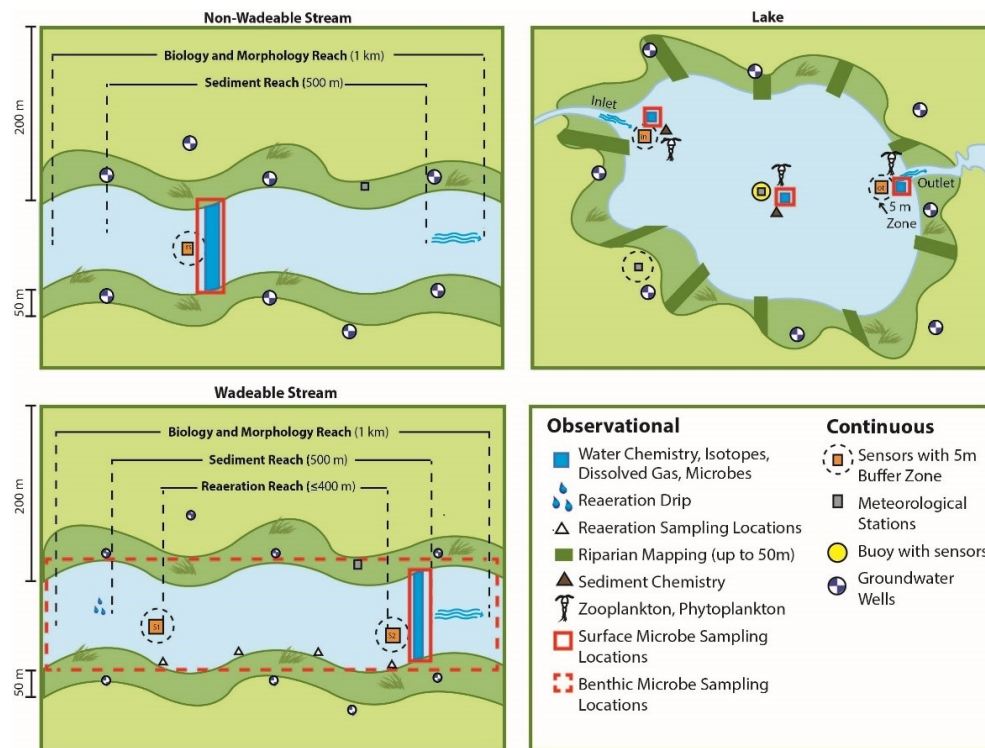


Figure 3: Generic NEON aquatic site layouts with microbial sampling locations highlighted in red boxes.

### 3.2 Temporal Sampling Design

At most terrestrial sites, soil group abundance sampling occurs 3 times per year in conjunction with the soil physical properties data product (DP1.10086). Two sampling bouts occur during periods of seasonal transitions (e.g. winter-spring or wet-dry), and once during the period of peak greenness (as measured by remote sensing data). At sites with short growing seasons (e.g. tundra and taiga), sampling occurs once annually during peak greenness.

Once every five years, a ‘coordinated’ bout occurs in which additional biogeochemical and isotopic measurements are made (DP1.10078), along with measurements of microbe biomass (DP1.10104) and nitrogen transformation rates (DP1.10080). During a coordinated bout, up to 2 soil horizons are sampled for microbial analyses to a maximum depth of 30 cm.

Surface water samples are collected monthly in wadeable streams, and every other month in lakes and rivers in conjunction with surface water chemistry sampling. Benthic microbe samples are collected three times per year, roughly spring, summer, and autumn at the same time as algal periphyton samples.

For all samples, the temporal resolution is that of a single collection date. For a comprehensive description of field methods, refer to TOS Protocol and Procedure: Soil Biogeochemical and Microbial Sampling (AD[10]) or AOS Protocol and Procedure: Aquatic Microbial Sampling (AD[12]) for soil and aquatic sampling protocols, respectively.

<p>Title: NEON User Guide to Microbe Group Abundances (NEON.DP1.10109; NEON.DP1.20277; NEON.DP1.20278)</p>	<p>Date: 02/27/2018</p>
<p>Author: Lee Stanish</p>	<p>Revision: A</p>

Descriptions of the upstream field data products for soil (DP1.10086), and aquatic surface water (DP1.20138) can be found in those respective Data Product User Guides.

### 3.3 Variables Reported

All variables reported from the field or laboratory technician (LO data) are listed in the file, NEON Raw Data Validation for Microbe Group Abundances (NEON.DP0.10109) (AD[04]). All variables reported in the published data (L1 data) are also provided separately in the following files:

- NEON Data Variables for Soil Microbe Group Abundances (NEON.DP1.10109) (AD[05]).
- NEON Data Variables for Benthic Microbe Group Abundances (NEON.DP1.20277) (AD[06]).
- NEON Data Variables for Surface Water Microbe Group Abundances (NEON.DP1.20278) (AD[07]).

Field names have been standardized with Darwin Core terms (<http://rs.tdwg.org/dwc/>; accessed 16 February 2014), the Global Biodiversity Information Facility vocabularies (<http://rs.gbif.org/vocabulary/gbif/>; accessed 16 February 2014), the VegCore data dictionary (<https://projects.nceas.ucsb.edu/nceas/projects/bien/wiki/VegCore>; accessed 16 February 2014), where applicable. NEON TOS spatial data employs the World Geodetic System 1984 (WGS84) for its fundamental reference datum and GEOID09 for its reference gravitational ellipsoid. NEON Aquatic spatial data uses the 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 Spatial Resolution and Extent

The finest resolution at which spatial data are reported is a single sampling location. For soils, this corresponds to a single X,Y coordinate location within a plot. For aquatics, this corresponds to a single station or habitat unit within a site.

#### 3.4.1 Soils

**sampleID** (unique ID given to the individual soil sampling location and horizon) → **plotID** (ID of plot within site) → **siteID** (ID of NEON site) → **domainID** (ID of a NEON domain).

The spatial data are located in the data product Soil Physical Properties, distributed periodic (DP1.10086), in the table *sls\_soilCoreCollection*. The spatial data are measured at the plot *centroid*, however, a more precise measurement may be desired by calculating the offset from the plot centroid using the variables **coreCoordinateX** and **coreCoordinateY**. Refer to the User Guide for Soil Physical Properties, distributed periodic, for more information and instructions.

### 3.4.2 Aquatics

**namedLocation** (unique ID given to the location within a site) → **siteID** (ID of NEON site) → **domainID** (ID of a NEON domain).

The spatial data can be found in the following Data Products:

- Surface water samples: Surface water microbe cell count (NEON.DP1.20138), in the table **amc\_fieldSuperParent**.
- Benthic samples: Benthic microbe marker gene sequences (NEON.DP1.20086), in the field data table **amb\_fieldParent**.

### 3.5 Temporal Resolution and Extent

The finest resolution at which temporal data are reported is the **collectDate**, the date and time of day when the sample was 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>

### 3.6 Associated Data Streams

This section describes the data products that are directly linked or closely related to the microbe group abundances data products.

#### 3.6.1 Soils

Soil data are derived from subsamples collected during soil biogeochemical and microbial sampling and include numerous related data products:

- Soil Physical Properties, distributed periodic (DP1.10086) - Field data, including soil moisture and pH, associated with a soil sample. These data are linked to the group abundances data by the **geneticSampleID**.
- Soil microbe community composition (NEON.DP1.10081) - Microbial community composition data derived from marker gene sequencing. The **dnaSampleID** variable in the tables **mcc\_soilTaxonTable\_16S** and **mcc\_soilTaxonTable\_ITS** may be used to link data in this product to soil microbe group abundances data.
- Soil microbe marker gene sequences (NEON.DP1.10108): Microbial 16S and ITS sequence data. The **dnaSampleID** variable in the tables **mmg\_soilDnaExtraction**, **mmg\_soilPcrAmplification** and **mmg\_soilMarkerGeneSequencing** can be used to link data in this product to the soil microbe group abundances data.
- Soil microbe biomass (NEON.DP1.10104) - Microbial biomass as measured by PLFA. Use information in the Soil Physical Properties data product (NEON.DP1.10086, table **sls\_soilCoreCollection**) to obtain the **biomassID** corresponding to the **sampleID**. The **sampleID** will map to a corresponding **geneticSampleID**, which can then be used to link data in the two data products.

- Soil inorganic nitrogen pools and transformations (NEON.DP1.10080) - Measurements derived by field incubations of soil cores or buried bags. As described for soil microbe biomass, use the **sampleID** from table **sls\_soilCoreCollection** to link these data products.
- Soil chemical properties (Distributed periodic) (NEON.DP1.10078) - Measurements of soil carbon and nitrogen. As with soil microbe biomass, the corresponding **sampleID** can be used to link data.
- Soil stable isotopes (Distributed periodic) (NEON.DP1.10100) - Measurements of soil carbon and nitrogen stable isotopes. As with soil microbe biomass, the corresponding **sampleID** can be used to link data.

### 3.6.2 Aquatics

Aquatic data are derived from samples collected in conjunction with other physical, chemical, and biological measurements. These include:

- Surface water microbes field data are found in the Aquatic Cell Counts data product (NEON.DP1.20138). The field **geneticSampleID** within the table **amc\_fieldCellCounts** can be used to link these data products.
- Benthic microbes field data are part of the download package for the Benthic microbe marker gene sequences data product (NEON.DP1.20280), and can be linked by the **geneticSampleID**.
- Chemical properties of surface water (NEON.DP1.20093) - Measurements of chemical constituents in water. The field **parentSampleID** in the table **swc\_fieldSuperParent** can be used to link these data products.
- Periphyton, seston and phytoplankton collection (NEON.DP1.20166) - Field data associated with sample collection. The field **parentSampleID** in the table **alg\_fieldData** links to the **sampleID** in the table **amb\_fieldParent**, which can then be linked to this data product by the **geneticSampleID**.
- Periphyton, seston and phytoplankton chemical properties (NEON.DP1.20163): Measurements of chemical constituents of algal samples. The field **parentSampleID** in the table **alg\_domainLabChemistry** links to the **sampleID** in the table **amb\_fieldParent**, which can then be linked to this data product by the **geneticSampleID**.
- Benthic (NEON.DP1.20086) and surface water (NEON.DP1.20141) microbe community composition: Taxonomic data derived from 16S and ITS marker gene sequencing. The field **dnaSampleID** in the tables **mcc\_benthicTaxonTable\_16S**, **mcc\_benthicTaxonTable\_ITS**, **mcc\_swTaxonTable\_16S** and **mcc\_swTaxonTable\_ITS** can be used to link these data to this data product.
- Benthic (NEON.DP1.20280) microbe marker gene sequence data. The field **geneticSampleID** in the tables **amb\_fieldParent** and **mmg\_benthicDnaExtraction** can be used to link these data to this data product.
- Surface water (NEON.DP1.20282) microbe marker gene sequences data. The field **geneticSampleID** in the tables **mmg\_swDnaExtraction** can be used to link these data to this data product.

### 3.7 Product Instances

For soil samples, a maximum of 10 plots will be sampled at every site one to three times per year. Most years, the surface soil horizon (organic or mineral) will be collected, while once every 5 years during a coordinated microbes/biogeochemistry bout, up to 2 soil horizons will be collected as separate samples. For each soil horizon sampled, 3 unique locations are collected at each plot, for up to 6 samples per plot. Thus, there will be 30-120 product instances generated per site per year.

Title: NEON User Guide to Microbe Group Abundances (NEON.DP1.10109; NEON.DP1.20277; NEON.DP1.20278)	Date: 02/27/2018
Author: Lee Stanish	Revision: A

Aquatic samples are collected at all aquatic NEON sites. For surface water sampling, a maximum of 4 sample locations will be sampled at every site 6-12 times per year, for a maximum of 24 product instances collected per site per year in a lake, and 12 product instances per site per year in a wadeable stream or river. Benthic microbial sampling occurs only at wadeable stream sites, where up to 8 samples are collected three times per year, for a maximum of 24 product instances per site per year.

### 3.8 Data Relationships

#### 3.8.1 Soils

The protocol dictates that each X,Y location sampled yields a unique **sampleID** per horizon per collectDate (day of year, local time) in the table *sls\_soilCoreCollection* for the data product Soil Physical Properties, distributed periodic (NEON.DP1.10086). Every bout type that includes microbes (e.g. the variable **boutType** includes the string 'microbe') should sample for group abundance analysis. A record from *sls\_soilCoreCollection* may have zero or one child records in table *mga\_soilGroupAbundances* of this data product.

##### Soil Physical Properties (NEON DP1.10086)

*sls\_soilCoreCollection.csv* -> One record expected per **sampleID**. Generates samples used in Soil microbe group abundances (NEON.DP1.10109), Soil microbe marker gene sequences (NEON.DP1.10108), Soil microbe community composition (NEON.DP1.10081), and Soil microbe biomass (NEON.DP1.10104). Additionally, subsamples generated from soil sampleIDs are used in Soil inorganic nitrogen pools and transformations (NEON.DP1.10080). Each **geneticSampleID** is a subsample of the parent **sampleID** and is sent for DNA extraction.

##### Soil Microbe Marker Gene Sequences (NEON.DP1.10108)

*mmg\_soilDnaExtraction.csv* -> This table contains the DNA extraction laboratory data. Data are linked by the **geneticSampleID**. There are one or more **dnaSampleIDs** expected per **geneticSampleID**, depending on the number of DNA extractions that occur on a sample provided to the lab. Duplicate records for an individual **dnaSampleID** should not exist. *Important Note:* This DNA extraction table is generic: samples that may not be relevant to the soil data product may appear in the data table. To limit the DNA extraction dataset to those that are relevant to the group abundances samples, filter the records in the *mmg\_soilDnaExtraction* table to include only those with a **dnaSampleID** that is also contained in the *mga\_soilGroupAbundances* table.

##### Soil Microbe Group Abundances (NEON.DP1.10109)

*mga\_soilGroupAbundances.csv* -> This table includes the gene copy number data for each sample. One record is expected per **dnaSampleID** per **targetTaxonGroup**.

*mga\_batchResults.csv* -> This table describes the batch-level data associated with a qPCR run. One record is expected per batch of samples analyzed (**batchID**), and links to the table *mga\_soilGroupAbundances* by the **batchID**. *Important Note:* The batch results table is generic for all soil and aquatic data: samples that may not be relevant to this data product may appear in the data table. To limit the dataset to those that are relevant to the soil group abundances data, filter the records to only those with **batchID**'s matching the **batchID**'s in the *mga\_soilGroupAbundances* table.

Title: NEON User Guide to Microbe Group Abundances (NEON.DP1.10109; NEON.DP1.20277; NEON.DP1.20278)	Date: 02/27/2018
Author: Lee Stanish	Revision: A

*mga\_labSummary.csv* -> This table describes the laboratory methods used to analyze samples, with **labSpecificStartDate** and **labSpecificEndDate** indicating the date range over which the methods apply. One record is expected per unique set of methods. The start and end dates can be used to filter the data in *mga\_soilGroupAbundances* using the fields **laboratoryName**, **processedDate**, and **targetTaxonGroup**.

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.*

### 3.8.2 Aquatics

#### 3.8.2.1 Surface Water

The protocol dictates that each namedLocation sampled yields a unique **parentSampleID**, one sample per collectDate (day of year, local time) in Surface water microbe cell count (DP1.20138), in the table *amc\_fieldSuperParent*. Each **parentSampleID** may be subsampled into one **geneticSampleID** that is used for microbial analyses, and an archive sample, described in the table *amc\_fieldCellCounts* within the Surface water microbe cell count product. These **geneticSampleIDs** are sent for DNA extraction such that the **geneticSampleID** from *amc\_fieldCellCounts* = **geneticSampleID** in *mmg\_swDnaExtraction*.

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.*

---

#### Surface Water Microbe Cell Count (NEON.DP1.20138)

*amc\_fieldSuperParent.csv* -> One record expected per namedLocation sampled and collectDate (day of year, local time), generates a unique **parentSampleID**.

*amc\_fieldCellCounts.csv* -> One record expected per namedLocation per collectDate (day of year, local time). Record represents a subsample (**geneticSampleID**) of the field-collected samples (**parentSampleID**). Depending on the time of year, each record generates zero or one **geneticSampleIDs**, corresponding to the Surface water microbe metagenome sequences (NEON.DP1.10107) variable **geneticSampleID** in the table *mmg\_swDnaExtraction*.

#### Surface Water Microbe Marker Gene Sequences (NEON.DP1.20282)

*mmg\_swDnaExtraction.csv* -> This table contains the DNA extraction laboratory data. Data are linked by the **geneticSampleID**. There are one or more **dnaSampleIDs** expected per **geneticSampleID**, depending on the number of DNA extractions that occur on a sample provided to the lab. Duplicate records for an individual **dnaSampleID** should not exist.

#### Surface Water Microbe Group Abundances (NEON.DP1.20278)

*mga\_swGroupAbundances.csv* -> This table includes the gene copy number data for each sample. One record expected per **dnaSampleID** per **targetTaxonGroup**.

*mga\_batchResults.csv* -> This table describes the batch-level data associated with a qPCR run. One record is expected per batch of samples analyzed (**batchID**), and links to the table *mga\_swGroupAbundances* by the **batchID**. *Important Note:* The batch results table is generic for all soil and aquatic data: samples that may not

<p>Title: NEON User Guide to Microbe Group Abundances (NEON.DP1.10109; NEON.DP1.20277; NEON.DP1.20278)</p>	<p>Date: 02/27/2018</p>
<p>Author: Lee Stanish</p>	<p>Revision: A</p>

be relevant to this data product may appear in the data table. To limit the dataset to those that are relevant to the soil group abundances data, filter the records to only those with **batchID**'s matching the **batchID**'s in the **mga\_swGroupAbundances\_** table.

*mga\_labSummary.csv* -> This table describes the laboratory methods used to analyze samples, with **labSpecificStartDate** and **labSpecificEndDate** indicating the date range over which the methods apply. One record is expected per unique set of methods. The start and end dates can be used to filter the data in **mga\_swGroupAbundances** using the fields **laboratoryName**, **processedDate**, and **targetTaxonGroup**.

### 3.8.2.2 Benthic

The protocol dictates that each namedLocation sampled yields a unique **sampleID**, one sample per collectDate (day of year, local time) in Benthic microbe marker gene sequences (DP1.20280), in the table **amb\_fieldParent**. Each **sampleID** may be subsampled into one **geneticSampleID** that is used for microbial analyses, and an archive sample, described in the same table. These **geneticSampleID**s are sent for DNA extraction such that the **geneticSampleID** from **amb\_fieldParent** = **geneticSampleID** in **mmg\_benthicDnaExtraction**.

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.*

---

#### Benthic Microbe Marker Gene Sequences (NEON.DP1.20280)

*amb\_fieldParent.csv* -> One record expected per namedLocation and collectDate (day of year, local time), and generates a unique **sampleID**. Record represents a subsample (**geneticSampleID**) of the field-collected sample.

*mmg\_benthicDnaExtraction.csv* -> This table contains the DNA extraction laboratory data. Data are linked by the **geneticSampleID**. There are one or more **dnaSampleID**s expected per **geneticSampleID**, depending on the number of DNA extractions that occur on a sample provided to the lab. Duplicate records for an individual **dnaSampleID** should not exist.

#### Benthic Microbe Group Abundances (NEON.DP1.20277)

*mga\_benthicGroupAbundances.csv* -> This table includes the gene copy number data for each sample. One record expected per **dnaSampleID** per **targetTaxonGroup**.

*mga\_batchResults.csv* -> This table describes the batch-level data associated with a qPCR run. One record is expected per batch of samples analyzed (**batchID**), and links to the table **mga\_benthicGroupAbundances** by the **batchID**. *Important Note:* The batch results table is generic for all soil and aquatic data: samples that may not be relevant to this data product may appear in the data table. To limit the dataset to those that are relevant to the soil group abundances data, filter the records to only those with **batchID**'s matching the **batchID**'s in the **mga\_benthicGroupAbundances\_** table.

*mga\_labSummary.csv* -> This table describes the laboratory methods used to analyze samples, with **labSpecificStartDate** and **labSpecificEndDate** indicating the date range over which the methods apply. One record is expected per unique set of methods. The start and end dates can be used to filter the data in **mga\_swGroupAbundances** using the fields **laboratoryName**, **processedDate**, and **targetTaxonGroup**.

<p>Title: NEON User Guide to Microbe Group Abundances (NEON.DP1.10109; NEON.DP1.20277; NEON.DP1.20278)</p>	<p>Date: 02/27/2018</p>
<p>Author: Lee Stanish</p>	<p>Revision: A</p>

## 4 DATA QUALITY

### 4.1 Data Entry Constraint and Validation

Constraints and data validation 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 document NEON Raw Data Validation for Microbe Group Abundances (NEON.DP0.10109), provided with every download of this data product. Contained within this file is a field named 'entryValidationRulesParser', which describes syntactically the validation rules for each field built into the data ingest validation. 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[16]).

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

To the extent possible, microbe group abundances data follow the *Minimum Information of Quantitative Real-Time PCR Experiments (MIQE)* metadata and QA/QC reporting guidelines outlined in Bustin et al. (2009). The MIQE standards define the essential and desired metadata to be reported for a qPCR reaction and include parameters related to experimental design, target gene information, oligonucleotides, the SOP/protocol, qPCR data validation and data analysis.

For each data product (soil, surface water, and benthic), the data table *mga\_groupAbundances* presents the abundance data for each sample as a unique record for each targetTaxonGroup. For example, samples that are analyzed for **targetTaxonGroups** 'bacteria', 'archaea' and 'fungi' separately should contain 3 records, while samples that are analyzed for the **targetTaxonGroups** 'bacteria and archaea' and 'fungi' should contain 2 records.

Following laboratory submission of metadata into the NEON automated data ingest process, 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[15]).

### 4.3 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.



<i>Title:</i> NEON User Guide to Microbe Group Abundances (NEON.DP1.10109; NEON.DP1.20277; NEON.DP1.20278)	<i>Date:</i> 02/27/2018
<i>Author:</i> Lee Stanish	<i>Revision:</i> A

#### 4.4 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

#### 4.5 Analytical Facility Data Quality

Data analyses conducted on qPCR data conform to the current data quality standards used by practitioners. The data table *mga\_labSummary* (available in the expanded package for each data product download) provides the general analytical approach as well as linkages to the SOP's and related documentation on the long-term methods used during the specified period of time.

## 5 REFERENCES

1. Bustin, S. A., V. Benes, J. A. Garson, J. Hellemans, J. Huggett, M. Kubista, R. Mueller, et al. 2009. The MIQE Guidelines: Minimum Information for Publication of Quantitative Real-Time PCR Experiments. *Clinical Chemistry* 55:611-622. <https://doi.org/10.1373/clinchem.2008.112797>.
2. Ginzinger, D. G. 2002. Gene quantification using real-time quantitative PCR: An emerging technology hits the mainstream. *Experimental Hematology* 30:503-512.