

- Data were retrieved that were designated as either 'grab' samples and 'composite' samples (mean result only).
  - No limits were specified for sample depths.
  - Data were retrieved for all fifty states, Puerto Rico, and the District of Columbia.
  - The time period specified for data retrieval was January 1990 to September 1998.
  - No data marked as "Retired Data" (i.e., data from a generally unknown source) were retrieved.
  - Data marked as "National Urban Runoff data" (i.e., data associated with sampling conducted after storm events to assess nonpoint source pollutants) were included in the retrieval. Such data are part of STORET's 'Archived' data.
  - Intensive survey data (i.e., data collected as part of specific studies) were retrieved.
2. Any values falling below the 1st percentile and any values falling above the 99th percentile were transformed into 'missing' values (i.e., values were effectively removed from the data set, but were not permanently eliminated).
  3. Based on the STORET 'Remark Code' associated with each retrieved data point, the following rules were applied (Table 2):

<b>TABLE 2: STORET REMARK CODE RULES</b>	
<b>STORET Remark Code</b>	<b>Keep or Delete Data Point</b>
blank - Data not remarked.	Keep
A - Value reported is the mean of two or more determinations.	Keep
B - Results based upon colony counts outside the acceptable ranges.	Delete
C - Calculated. Value stored was not measured directly, but was calculated from other data available.	Keep
D - Field measurement.	Keep

E - Extra sample taken in compositing process.	Delete
F - In the case of species, F indicates female sex.	Delete
G - Value reported is the maximum of two or more determinations.	Delete
H - Value based on field kit determination; results may not be accurate.	Delete
I - The value reported is less than the practical quantification limit and greater than or equal to the method detection limit.	Keep, but used one-half the reported value as the new value.
J - Estimated. Value shown is not a result of analytical measurement.	Delete
K - Off-scale low. Actual value not known, but known to be less than value shown.	Keep, but used one-half the reported value as the new value.
L - Off-scale high. Actual value not known, but known to be greater than value shown.	Keep
M - Presence of material verified, but not quantified. Indicates a positive detection, at a level too low to permit accurate quantification.	Keep, but used one half the reported value as the new value.
N - Presumptive evidence of presence of material.	Delete
O - Sample for, but analysis lost. Accompanying value is not meaningful for analysis.	Delete
P - Too numerous to count.	Delete
Q - Sample held beyond normal holding time.	Delete
R - Significant rain in the past 48 hours.	Delete
S - Laboratory test.	Keep
T - Value reported is less than the criteria of detection.	Keep, but replaced reported value with 0.

U - Material was analyzed for, but not detected. Value stored is the limit of detection for the process in use.	Keep, but replaced reported value with 0.
V - Indicates the analyte was detected in both the sample and associated method blank.	Delete
W - Value observed is less than the lowest value reportable under remark "T."	Keep, but replaced reported value with 0.
X - Value is quasi vertically-integrated sample.	No data point with this remark code in data set.
Y - Laboratory analysis from unpreserved sample. Data may not be accurate.	Delete
Z - Too many colonies were present to count.	Delete
<p>If a parameter (excluding water temperature) value was less than or equal to zero and no remark code was present, the value was transformed into a missing value.  Rationale - Parameter concentrations should never be zero without a proper explanation. A method detection limit should at least be listed.</p>	

4. Station records were eliminated from the data set if any of the following descriptors were present within the "Station Type" parameter:

- ▶ **MONITR** - Source monitoring site, which monitors a known problem or to detect a specific problem.
- ▶ **HAZARD** - Site of hazardous or toxic wastes or substances.
- ▶ **ANPOOL** - Anchialine pool, underground pools with subsurface connections to watertable and ocean.
- ▶ **DOWN** - Downstream (i.e., within a potentially polluted area) from a facility which has a potential to pollute.
- ▶ **IMPDMT** - Impoundment. Includes waste pits, treatment lagoons, and settling and evaporation ponds.
- ▶ **STMSWR** - Storm water sewer.
- ▶ **LNDFL** - Landfill.
- ▶ **CMBMI** - Combined municipal and industrial facilities.
- ▶ **CMBSRC** - Combined source (intake and outfall).

Rationale - these descriptors potentially indicate a station location that at which an ambient water sample would not be obtained (i.e., such sampling locations are potentially

biased) or the sample location is not located within one of the designated water body types (i.e., ANPOOL).

5. Station records were eliminated from data set if the station location did not fall within any established cataloging unit boundaries based on their latitude and longitude.
6. Using nutrient ecoregion GIS coverage provided by USEPA, all station locations with latitude and longitude coordinates were tagged with a nutrient ecoregion identifier (nutrient region identifiers are values 1 - 14) and the associated nutrient ecoregion name. Because no nutrient ecoregions exist for Alaska, Hawaii, and Puerto Rico, stations located in these states were tagged with "dummy" nutrient ecoregion numbers (20 = Alaska, 21 = Hawaii, 22 = Puerto Rico).
7. Using information provided by TVA, 59 station locations that were marked as 'stream' locations under the "Station Type" parameter were changed to 'reservoir' locations.
8. The nutrient data retrieved from STORET were assessed for the presence of duplicate data records. The duplicate data identification process consisted of three steps: 1) identification of records that matched exactly in terms of each variable retrieved; 2) identification of records that matched exactly in terms of each variable retrieved except for their station identification numbers; and 3) identification of records that matched exactly in terms of each variable retrieved except for their collecting agency codes. The data duplication assessment procedures were conducted using SAS programs. Prior to initiating the data duplication assessment process, the STORET nutrient data set contained:

41,210 station records  
924,420 sample records

- Identification of exactly matching records  
All data records were sorted to identify those records that matched exactly. For two records to match exactly, all variables retrieved had to be the same. For example, they had to have the same water quality parameters, parameter results and associated remark codes, and have the same station data item and sample data item information. Exactly matching records were considered to be exact duplicates, and one duplicate record of each identified matching set were eliminated from the nutrient data set. A total of 924 sample records identified as duplicates by this process were eliminated from the data set.
- Identification of matching records with the exception of station identification number  
All data records were sorted to identify those records that matched exactly except for their station identification number (i.e., they had the same water quality

parameters, parameter results and associated remark codes, and the same station and sample data item information with the exception of station identification number). Although the station identification numbers were different, the latitude and longitude for the stations were the same indicating a duplication of station data due to the existence of two station identification numbers for the same station. For each set of matching records, one of the station identification numbers was randomly selected and its associated data were eliminated from the data set. A total of 686 sample records were eliminated from the data set through this process.

- Identification of matching records with the exception of collecting agency codes  
All data records were sorted to identify those records that matched exactly except for their collecting agency codes (i.e., they had the same water quality parameters, parameter results and associated remark codes, and the same station and sample data item information with the exception of agency code). The presence of two matching data records each with a different agency code attached to it suggested that one agency had utilized data collected by the other agency and had entered the data into STORET without realizing that it already had been placed in STORET by the other agency. No matching records with greater than two different agency codes were identified. For determining which record to delete from the data set, the following rules were developed:
  - ▶ If one of the matching records had a USGS agency code, the USGS record was retained and the other record was deleted.
  - ▶ Higher level agency monitoring program data were retained. For example, federal program data (indicated by a "1" at the beginning of the STORET agency code) were retained against state (indicated by a "2") and local (indicated by values higher than 2) program data.
  - ▶ If two matching records had the same level agency code, the record from the agency with the greater number of overall observations (potentially indicating the data set as the source data set) was retained.

A total of 2,915 sample records were eliminated through this process.

As a result of the duplicate data identification process, a total of 4,525 sample records and 36 individual station records were removed from the STORET nutrient data set. The resulting nutrient data set contains the following:

41,174 station records  
919,895 sample records

## **APPENDIX B**

### **Process for Adding Aggregate Nutrient Ecoregions and Level III Ecoregions**

Steps for assigning Level III ecoregions and aggregate nutrient ecoregion codes and names to the Nutrient Criteria Database (performed using ESRI's ARCView v 3.2 and its GeoProcessing Wizard). This process is performed twice; once for the Level III ecoregions and once for the aggregate nutrient ecoregions:

- Add the station .dbf data table, with latitude and longitude data, to project by 'Add Event Theme'
- Convert to the shapefile format
- Create 'stcojoin' field, populate the 'stcojoin' field with the following formula: 'County.LCase+State.LCase'
- Add field 'stco\_flag' to the station shapefile
- Spatially join the station data with the county shapefile (cntys\_jned.shp)
- Select 'stcojoin' (station shapefile) field = 'stco\_join2' (county shapefile) field
- Calculate stco\_flag = 0 for selected features
- Step through all blank stco\_flag records, assign the appropriate stco\_flags, see list on the following page
- Select all stco\_flags = 4 or 7, switch selection
- Calculate ctyfips (station) to cntyfips (county)
- Stop editing and save edits, remove all joins
- Add in 2 new fields 'x-coord1' and 'y-coord1' into station table
- Select all stco\_flags = 1, 2, and 6
- Link county coverage with station coverage
- Populate 'x-coord1' and 'y-coord1' with 'x-coord' and 'y-coord' from county coverage
- Select all stco\_flags = 1, 2, and 6, export to new .dbf file
- Add new .dbf file as event theme
- Convert to shapefile format
- Add the following fields to both tables (original station and station126 shapefiles): 'eco\_omer', 'name\_omer', 'dis\_aggr', 'code\_aggr', 'name\_aggr'
- Spatially join station126 and eco-omer coverage
- Populate the 'eco\_omer' field with the 'eco' value
- Repeat the previous step using the nearest method (line coverage) to determine ecoregion assignment for the line coverage, if some records are blank
- Spatially join the ecoregion line coverage to station coverage, link the LPoly# (from the spatially joined table) to Poly# (of the ecoregion polygon coverage)
- Populate the Eco fields with the appropriate information.
- Follow the same steps to the Rpoly#
- Remove all table joins
- Link the usco-om table with station126 table and populate 'name-omer' field
- Spatially join station aggr coverage and populate the rest of the fields. Follow the same procedures as outlined above
- Remove all joins
- Make sure the new Eco field added into the station126 shapefile are different than

- the ones in the original station shapefile
- Join station126 and station coverage by station-id
- Populate all the Eco fields in the original station coverage
- Remove all joins
- Save table
- Make sure that all ctyfips records are populated; the county shapefile may have to be joined to populate the records, if the stco\_flag = 4
- Create 2 new fields, 'NewCounty' and 'NewState'
- Populate these new fields with a spatial join to the county coverage
- Select by feature (ecoregion shapefile) all of the records in the station shapefile
- Switch selection (to get records outside of the ecoregion shapefile)
- If any of the selected records have stco\_flag = 0 (they are outside the ecoregion shapefile boundary), calculate them to stco\_flag = 3

### **stco\_flags (state/county flags in order of importance)**

- 0 The state and county values from the data set matched the state and county values from the spatial join.  
(Ecoregions were assigned based on the latitude/longitude coordinates.)
- 1 The state and county values from the data set did not match the state and county values from the spatial join, but the point was inside the county coverage boundary.  
(Ecoregions were assigned based on the county centroid.)
- 2 The state and county values from the data set did not match the state and county values from the spatial join because the point was outside the county coverage boundary; therefore, there was nothing to compare to the point (i.e., the point falls in the ocean/Canada/Mexico). This occurred for some coastal samples.  
(Ecoregions were assigned based on the county centroid.)
- 3 The state and county values from the data set matched the state and county from the spatial join, but the point was outside the ecoregion boundary.  
(Ecoregions were assigned to the closest ecoregion to the point.)  
(No ecoregions were assigned to AK, HI, PR, BC, and GU.)
- 4 Latitude/longitude coordinates were provided, but there was no county information.  
(Ecoregions were assigned based on the latitude/longitude coordinates.)
- 5 The state and county values from the data set did not match the state and county values from the spatial join due to spelling or naming convention errors. The matches were performed manually.  
(Ecoregions were assigned based on the latitude/longitude coordinates.)
- 6 No latitude/longitude coordinates were provided, only state and county information was available.  
(Ecoregions were assigned based on the county centroid.)
- 7 No latitude/longitude coordinates were provided, only state information was available; therefore, no matches were possible.  
(Ecoregions were not assigned. Data is not included in the analysis.)

## **APPENDIX C**

### **Glossary**

Coefficient of Variation- Equal to the standard deviation divided by the mean multiplied by 100.

Maximum- The highest value.

Mean- The arithmetic average.

Median- The 50th percentile or middle value. Half of the values are above the median, and half of the values are below the median.

Minimum- The lowest value.

Standard Deviation- Equal to the square root of the variance with the variance defined as the sum of the squared deviations divided by the sample size minus one.

Standard Error- Standard error of the mean is equal to the standard deviation divided by the square root of the sample size.