

# 5. Bd-Capture Workflow

Cob Staines

2026-02-18

## Table of contents

<b>Motivation</b>	<b>1</b>
<b>R</b>	<b>2</b>
Setup . . . . .	2
Data discovery and pulling . . . . .	2
Point to support tables . . . . .	2
Join all data of interest . . . . .	3
Select columns of interest, filter to date . . . . .	5
Explore # of filtered observations by life stage, then filter again . . . . .	6
Pull data . . . . .	6
Disconnect . . . . .	10
<b>Python</b>	<b>10</b>
Setup . . . . .	10
Data discovery and pulling . . . . .	11
Point to support tables . . . . .	11
Join all data of interest . . . . .	11
Filter to date, select columns of interest . . . . .	12
Explore # of filtered observations by life stage, then filter again . . . . .	13
Pull data . . . . .	14
Disconnect . . . . .	14

*This tutorial is available as a [.qmd on Github](#).*

## Motivation

- Run through a complete workflow including: database connection, data discovery, data pulling, and data manipulation
- Demonstrate how to connect Capture data with sample data (in this case Bd qPCR results)

# R

Let's run through a realistic data scenario from the beginning.

Suppose we are interested in Capture and Bd qPCR data. Specifically, we want to compare Bd qPCR results for juvenile-to-adult individuals, captured in 2015 or later, across species and sites.

## Setup

```
# minimal packages for RIBBiTR DB data discovery
librarian::shelf(tidyverse, dbplyr, RPostgres, DBI, RIBBiTR-BII/ribbitrrr)

# establish database connection
dbcon <- hopToDB("ribbitr")
```

Connecting to 'ribbitr'... Success!

```
# load table metadata
mdt <- tbl(dbcon, Id("public", "all_tables")) %>%
  filter(table_schema == "survey_data") %>%
  collect()

# load column metadata
mdc <- tbl(dbcon, Id("public", "all_columns")) %>%
  filter(table_schema == "survey_data") %>%
  collect()
```

## Data discovery and pulling

Looking at the [survey\\_data schema diagram](#), we can browse to see which tables and columns we want. We can also consult the table or column metadata. The two observation tables with the primary data of interest are called “capture” and “bd\_qpcr\_results”.

## Point to support tables

```
# pointers for all tables of interest
db_bdqpcr = tbl(dbcon, Id("survey_data", "bd_qpcr_results"))
db_sample = tbl(dbcon, Id("survey_data", "sample"))
db_capture = tbl(dbcon, Id("survey_data", "capture"))
```

```

db_survey = tbl(dbcon, Id("survey_data", "survey"))
db_visit = tbl(dbcon, Id("survey_data", "visit"))
db_site = tbl(dbcon, Id("survey_data", "site"))
db_region = tbl(dbcon, Id("survey_data", "region"))
db_country = tbl(dbcon, Id("survey_data", "country"))

```

## Join all data of interest

In this case we only want to consider cases for which we have both capture *and* Bd qPCR data. An inner join across the Bd, sample, and capture tables will keep only values which are common between them. Of these data, we want to return all supporting info, so we will left join to the remaining support tables.

```

# inner join capture and bd samples
# left join supporting tables
data_bd_capture = db_bdqpcr %>%
  inner_join(db_sample, by = "sample_id") %>%
  inner_join(db_capture, by = "capture_id") %>%
  left_join(db_survey, by = "survey_id") %>%
  left_join(db_visit, by = "visit_id") %>%
  left_join(db_site, by = "site_id") %>%
  left_join(db_region, by = "region_id") %>%
  left_join(db_country, by = "country_id")

# see what columns are available
colnames(data_bd_capture)

```

[1] "result_id"	"sample_id"
[3] "sample_name_bd"	"bd_detected"
[5] "replicates"	"replicate_id"
[7] "bd_cycle_quant"	"bd_target_quant"
[9] "total_qpcr_volume_uL"	"qpcr_dilution_factor"
[11] "template_dna_volume_uL"	"extract_volume_uL"
[13] "bd_target_quant_per_swab"	"bd_its1_copies_per_swab"
[15] "swab_type"	"standard_target_type"
[17] "standard"	"master_mix"
[19] "extraction_plate_name"	"extraction_date"
[21] "extraction_kit"	"extraction_lab"
[23] "qpcr_plate_name"	"qpcr_well"
[25] "qpcr_date"	"qpcr_machine"
[27] "qpcr_lab"	"comments_qpcr"
[29] "its1_copies_per_standard_unit"	"bsal_detected"
[31] "sample_name"	"sample_type"

[33]	"capture_id"	"sample_name_conflict"
[35]	"negative_control_group_id"	"negative_control"
[37]	"taxon_capture"	"time_of_capture"
[39]	"capture_transect_m"	"microhabitat_type"
[41]	"body_temp_c"	"substrate_temp_c"
[43]	"svl_mm"	"body_mass_g"
[45]	"life_stage"	"sex"
[47]	"capture_animal_state"	"comments_capture"
[49]	"photo"	"photo_id"
[51]	"microhabitat_detailed"	"body_and_bag_mass_g"
[53]	"bag_mass_g"	"marked"
[55]	"capture_utme"	"capture_utmn"
[57]	"capture_type"	"observer_capture"
[59]	"bag_id"	"processor"
[61]	"cmr_id"	"microhabitat_notes"
[63]	"tail_length_mm"	"bucket"
[65]	"inside_outside_serdp"	"temp_gun"
[67]	"clearcut"	"number_of_mites"
[69]	"flir"	"tad_stage"
[71]	"survey_id"	"microhabitat_wet"
[73]	"capture_utm_zone"	"capture_latitude"
[75]	"capture_longitude"	"timestamp_of_capture_utc"
[77]	"treatment"	"start_time"
[79]	"end_time"	"detection_type"
[81]	"duration_minutes"	"observers_survey"
[83]	"comments_survey"	"description"
[85]	"survey_quality"	"transect"
[87]	"number_observers"	"visit_id"
[89]	"start_timestamp_utc"	"end_timestamp_utc"
[91]	"detection_subtype"	"date"
[93]	"time_of_day"	"campaign"
[95]	"visit_status"	"comments_visit"
[97]	"site_id"	"visit_lab"
[99]	"project_id"	"site"
[101]	"site_utm_zone"	"site_utme"
[103]	"site_utmn"	"area_sqr_m"
[105]	"site_code"	"site_elevation_m"
[107]	"depth_m"	"topo"
[109]	"wilderness"	"site_comments"
[111]	"region_id"	"site_name_alt"
[113]	"site_latitude"	"site_longitude"
[115]	"geographic_area"	"geographic_area_type"
[117]	"region"	"country_id"
[119]	"time_zone"	"country"
[121]	"iso_country_code"	

```
# we can also see which columns come from specified tables, for context  
colnames(db_bdqpcr)
```

```
[1] "result_id"                      "sample_id"  
[3] "sample_name_bd"                 "bd_detected"  
[5] "replicates"                    "replicate_id"  
[7] "bd_cycle_quant"                "bd_target_quant"  
[9] "total_qpcr_volume_uL"          "qpcr_dilution_factor"  
[11] "template_dna_volume_uL"        "extract_volume_uL"  
[13] "bd_target_quant_per_swab"      "bd_its1_copies_per_swab"  
[15] "swab_type"                     "standard_target_type"  
[17] "standard"                      "master_mix"  
[19] "extraction_plate_name"        "extraction_date"  
[21] "extraction_kit"               "extraction_lab"  
[23] "qpcr_plate_name"              "qpcr_well"  
[25] "qpcr_date"                     "qpcr_machine"  
[27] "qpcr_lab"                      "comments_qpcr"  
[29] "its1_copies_per_standard_unit" "bsal_detected"
```

### Select columns of interest, filter to date

```
# pull data from database  
data_bd_capture_2015 = data_bd_capture %>%  
  # filter to dates of interest  
  filter(date >= "2015-01-01") %>%  
  # select columns of interest  
  select(capture_id,  
         taxon_capture,  
         life_stage,  
         svl_mm,  
         body_mass_g,  
         survey_id,  
         cmr_id,  
         sample_id,  
         sample_name_bd,  
         bd_detected,  
         bd_cycle_quant,  
         bd_target_quant,  
         bd_target_quant_per_swab,  
         bd_its1_copies_per_swab,  
         comments_capture,  
         comments_qpcr,
```

```
date,  
site,  
region,  
country)
```

### Explore # of filtered observations by life stage, then filter again

```
data_bd_capture_2015 %>%  
  select(life_stage) %>%  
  group_by(life_stage) %>%  
  summarise(row_count = n()) %>%  
  arrange(desc(row_count)) %>%  
  collect()
```

```
# A tibble: 10 x 2  
  life_stage     row_count  
  <chr>           <int64>  
1 adult            34884  
2 juvenile          4397  
3 tadpole           2763  
4 <NA>              2584  
5 subadult          1862  
6 metamorph         429  
7 larva              417  
8 larvae             218  
9 metamorphosed      28  
10 eggmass            7
```

```
data_bd_capture_2015_life <- data_bd_capture_2015 %>%  
  filter(life_stage %in% c("juvenile",  
                           "subadult",  
                           "adult"),  
         !is.na(life_stage))
```

### Pull data

```
# inspect our SQL "shopping list"  
sql_render(data_bd_capture_2015_life)
```

```

<SQL> SELECT
    "capture_id",
    "taxon_capture",
    "life_stage",
    "svl_mm",
    "body_mass_g",
    "survey_id",
    "cmr_id",
    "sample_id",
    "sample_name_bd",
    "bd_detected",
    "bd_cycle_quant",
    "bd_target_quant",
    "bd_target_quant_per_swab",
    "bd_its1_copies_per_swab",
    "comments_capture",
    "comments_qpcr",
    "date",
    "site",
    "region",
    "country"
FROM (
    SELECT
        "bd_qpcr_results".*,
        "sample_name",
        "sample_type",
        "sample"."capture_id" AS "capture_id",
        "sample_name_conflict",
        "negative_control_group_id",
        "negative_control",
        "taxon_capture",
        "time_of_capture",
        "capture_transect_m",
        "microhabitat_type",
        "body_temp_c",
        "substrate_temp_c",
        "svl_mm",
        "body_mass_g",
        "life_stage",
        "sex",
        "capture_animal_state",
        "comments_capture",
        "photo",
        "photo_id",
        "microhabitat_detailed",

```

```
"body_and_bag_mass_g",
"bag_mass_g",
"marked",
"capture_utme",
"capture_utmn",
"capture_type",
"observer_capture",
"bag_id",
"processor",
"cmr_id",
"microhabitat_notes",
"tail_length_mm",
"bucket",
"inside_outside_serdp",
"temp_gun",
"clearcut",
"number_of_mites",
"flir",
"tad_stage",
"capture"."survey_id" AS "survey_id",
"microhabitat_wet",
"capture_utm_zone",
"capture_latitude",
"capture_longitude",
"timestamp_of_capture_utc",
"treatment",
"start_time",
"end_time",
"detection_type",
"duration_minutes",
"observers_survey",
"comments_survey",
"description",
"survey_quality",
"transect",
"number_observers",
"survey"."visit_id" AS "visit_id",
"start_timestamp_utc",
"end_timestamp_utc",
"detection_subtype",
"date",
"time_of_day",
"campaign",
"visit_status",
"comments_visit",
```

```

"visit"."site_id" AS "site_id",
"visit_lab",
"project_id",
"site",
"site_utm_zone",
"site_utme",
"site_utmn",
"area_sqr_m",
"site_code",
"site_elevation_m",
"depth_m",
"topo",
"wilderness",
"site_comments",
"site"."region_id" AS "region_id",
"site_name_alt",
"site_latitude",
"site_longitude",
"geographic_area",
"geographic_area_type",
"region",
"region"."country_id" AS "country_id",
"time_zone",
"country",
"iso_country_code"
FROM "survey_data"."bd_qpcr_results"
INNER JOIN "survey_data"."sample"
    ON ("bd_qpcr_results"."sample_id" = "sample"."sample_id")
INNER JOIN "survey_data"."capture"
    ON ("sample"."capture_id" = "capture"."capture_id")
LEFT JOIN "survey_data"."survey"
    ON ("capture"."survey_id" = "survey"."survey_id")
LEFT JOIN "survey_data"."visit"
    ON ("survey"."visit_id" = "visit"."visit_id")
LEFT JOIN "survey_data"."site"
    ON ("visit"."site_id" = "site"."site_id")
LEFT JOIN "survey_data"."region"
    ON ("site"."region_id" = "region"."region_id")
LEFT JOIN "survey_data"."country"
    ON ("region"."country_id" = "country"."country_id")
) AS "q01"
WHERE
    ("date" >= '2015-01-01') AND
    ("life_stage" IN ('juvenile', 'subadult', 'adult')) AND
    (NOT(("life_stage" IS NULL)))

```

```

# collect data
data_bd_capture_final <- data_bd_capture_2015_life %>%
  collect()

head(data_bd_capture_final)

# A tibble: 6 x 20
  capture_id      taxon_capture life_stage svl_mm body_mass_g survey_id cmr_id
  <chr>          <chr>        <chr>       <dbl>     <dbl>    <chr>    <chr>
1 6eeebe74d-5325-44~ espadaranana_p~ adult       25.5      0.7 473bfaba~ <NA>
2 96fba6ff-677d-4b~ colostethus_~ adult       27.3      2.1 c500c2d7~ <NA>
3 05aac450-7c9e-4a~ colostethus_~ adult       26.5      1.5 60b7f2d7~ <NA>
4 f1ab9744-86a7-45~ boana_platan~ adult      55.0      15.3 5ca54b85~ <NA>
5 707f2783-6de8-40~ colostethus_~ adult       27.3      2.8 c500c2d7~ <NA>
6 e68a063e-c259-4c~ colostethus_~ adult       26.8      2.7 c500c2d7~ <NA>
# i 13 more variables: sample_id <chr>, sample_name_bd <chr>,
#   bd_detected <lgl>, bd_cycle_quant <dbl>, bd_target_quant <dbl>,
#   bd_target_quant_per_swab <dbl>, bd_its1_copies_per_swab <dbl>,
#   comments_capture <chr>, comments_qpcr <chr>, date <date>, site <chr>,
#   region <chr>, country <chr>

```

These data are ready to be analyzed and visualized!

## Disconnect

```
dbDisconnect(dbcon)
```

## Python

Let's run through a realistic data scenario from the beginning.

Suppose we are interested in Capture and Bd qPCR data. Specifically, we want to compare Bd qPCR results for juvenile-to-adult individuals, captured in 2015 or later, across species and sites.

## Setup

```

# minimal packages for RIBBiTR DB Workflow
import ibis
from ibis import _
import pandas as pd
import dbconfig
import db_access as db

# establish database connection
dbcon = ibis.postgres.connect(**dbconfig.ribbitr)

# load table metadata
mdt = dbcon.table(database = "public", name = "all_tables").to_pandas()

# load column metadata
mdc = (
    dbcon.table(database="public", name="all_columns")
    .filter(_.table_schema == 'survey_data')
    .to_pandas()
)

```

## Data discovery and pulling

Looking at the [survey\\_data schema diagram](#), we can browse to see which tables and columns we want. We can also consult the table or column metadata. The two observation tables with the primary data of interest are called `capture` and `bd_qpcr_results`.

## Point to support tables

```

# Pointers for all tables
db_bdqpcr = dbcon.table('bd_qpcr_results', database='survey_data')
db_sample = dbcon.table('sample', database='survey_data')
db_capture = dbcon.table('capture', database='survey_data')
db_survey = dbcon.table('survey', database='survey_data')
db_visit = dbcon.table('visit', database='survey_data')
db_site = dbcon.table('site', database='survey_data')
db_region = dbcon.table('region', database='survey_data')
db_country = dbcon.table('country', database='survey_data')

```

## Join all data of interest

In this case we only want to consider cases for which we have both capture *and* Bd qPCR data. An inner join across the Bd, sample, and capture tables will keep only values which

are common between them. Of these data, we want to return all supporting info, so we will left join to the remaining support tables.

```
# Recursive joins
data_bd_capture = (
    db_bdqpcr
    .inner_join(db_sample, db_bdqpcr.sample_id == db_sample.sample_id)
    .inner_join(db_capture, db_sample.capture_id == db_capture.capture_id)
    .left_join(db_survey, db_capture.survey_id == db_survey.survey_id)
    .left_join(db_visit, db_survey.visit_id == db_visit.visit_id)
    .left_join(db_site, db_visit.site_id == db_site.site_id)
    .left_join(db_region, db_site.region_id == db_region.region_id)
    .left_join(db_country, db_region.country_id == db_country.country_id)
)

# see what columns are available
data_bd_capture.columns

['result_id', 'sample_id', 'sample_name_bd', 'bd_detected', 'replicates', 'replicate_id', 'b
# we can also see which columns come from specified tables, for context
db_bdqpcr.columns

['result_id', 'sample_id', 'sample_name_bd', 'bd_detected', 'replicates', 'replicate_id', 'b
```

### Filter to date, select columns of interest

```
# capture table, filter, select
data_bd_capture_2015 = (
    data_bd_capture
    .filter(_.date >= '2015-01-01')
    .select([
        "capture_id",
        "taxon_capture",
        "life_stage",
        "svl_mm",
        "body_mass_g",
        "survey_id",
        "cmr_id",
        "sample_id",
        "sample_name_bd",
        "bd_detected",
```

```

    "bd_cycle_quant",
    "bd_target_quant",
    "bd_target_quant_per_swab",
    "bd_its1_copies_per_swab",
    "comments_capture",
    "comments_qpcr",
    "date",
    "site",
    "region",
    "country"
  ])
)

```

### Explore # of filtered observations by life stage, then filter again

```

# count by life stage
life_stage_counts = (
    data_bd_capture_2015
    .group_by('life_stage')
    .aggregate(row_count=_.count())
    .order_by(_.row_count.desc())
    .to_pandas()
)

```

```
print(life_stage_counts)
```

	life_stage	row_count
0	adult	34884
1	juvenile	4397
2	tadpole	2763
3	None	2584
4	subadult	1862
5	metamorph	429
6	larva	417
7	larvae	218
8	metamorphosed	28
9	eggmass	7

```

# filter to desired life stages
data_bd_capture_2015_life = (
    data_bd_capture_2015

```

```
.filter(_.life_stage.isin(['juvenile','subadult', 'adult']) & _.life_stage.notnull())
)
```

## Pull data

```
# inspect our SQL "shopping list"
data_bd_capture_2015_life.compile()
```

```
'SELECT "t16"."capture_id", "t16"."taxon_capture", "t16"."life_stage", "t16"."svl_mm", "t16"
```

```
# pull data
data_bd_capture_final = data_bd_capture_2015_life.to_pandas()

# preview data
data_bd_capture_final.head()
```

```
capture_id    ... country
0  96fba6ff-677d-4b8c-84e7-a17763b18523  ... panama
1  6eeebe74d-5325-44a0-bc9e-e19c930e7f2e  ... panama
2  05aac450-7c9e-4acf-9e98-80aa97750cd6  ... panama
3  f1ab9744-86a7-45d7-9014-f4a6da9c857f  ... panama
4  707f2783-6de8-4048-aa3a-1bbc3c149710  ... panama
```

```
[5 rows x 20 columns]
```

These data are ready to be analyzed and visualized!

## Disconnect

```
# close connection
dbcon.disconnect()
```

[<- 4. Table Joins | 6. Microclimate Workflow ->](#)