

Fric et al. Re-analysis Code

Vaughn Shirey & Elise Larsen

Current version 11/19/2020; initiated 3/9/2020

*

Begin Analysis

This code chunk sets up the workspace and loads necessary packages. If phest is not already installed, remove comment from install line.

```
# Load Libraries
library(tidyverse)
library(ggplot2)
library(ggExtra)
library(gridExtra)
#library(devtools); install_github("willpearse/phest")
library(phest)
library(readxl)
```

*

Data Import and Formatting

data.csv file was downloaded from <https://doi.org/10.6084/m9.figshare.9946934>
(<https://doi.org/10.6084/m9.figshare.9946934>)
(https://figshare.com/articles/Phenology_responses_of_temperate_butterflies_-_Supplementary_data/9946934)
(https://figshare.com/articles/Phenology_responses_of_temperate_butterflies_-_Supplementary_data/9946934)

This csv file contains the occurrence data used in Fric et al. (2020), which they downloaded from gbif. The file includes separate data tables for each dataset, which have been concatenated into one file. These data tables have the same fields but are not formatted as a single data table; individual datasets were all written into one data file, including headers and row indices in each dataset. This first set of code reformats the data & writes formatted data files.

```

all.data <- readLines("fric_supplements/data.csv")

#identify header rows
all.header.rows<-grep("decimalLongitude", all.data)

#check headers for consistency
uniqueheaders<-unique(all.data[all.header.rows])

# 2 versions! -> Get row numbers for "header 1"
header.rows1<-grep(uniqueheaders[1], all.data)
#Get row numbers for "header 2"
header.rows2<-setdiff(all.header.rows, header.rows1)

#Create row identifiers:
#0 is a header row, 1 is format 1 data, 2 is format 2 data
j<-rep(0,length(all.data))
for (i in all.header.rows) {
  #set index to the next header if it's not the last header; otherwise set to end of datafile +
  1
  if(i<max(all.header.rows)) {
    next_index<-min(all.header.rows[all.header.rows>i])
  }else { next_index<-length(all.data)+1 }

  #for data between header rows, set row index
  j[(i+1):(next_index-1)]<-ifelse(i%in%header.rows1,1,2)
}

#need to add a row index to the header text for new data files
newheader1<-paste('"row.index\\"",', uniqueheaders[1], sep="")
newheader2<-paste('"row.index\\"",', uniqueheaders[2], sep="")

#write data file
formatteddatafile1<-file("data/fric_data_header_1.txt")
writeLines(c(newheader1,all.data[which(j==1)]), formatteddatafile1)
close(formatteddatafile1)

formatteddatafile2<-file("data/fric_data_header_2.txt")
writeLines(c(newheader2,all.data[which(j==2)]), formatteddatafile2)
close(formatteddatafile2)
rm(list=ls())

#read back in the formatted data
data1<-read_csv("data/fric_data_header_1.txt")

```

```
## Parsed with column specification:
## cols(
##   row.index = col_double(),
##   name = col_character(),
##   decimalLongitude = col_double(),
##   decimalLatitude = col_double(),
##   year = col_double(),
##   month = col_double(),
##   country = col_character(),
##   day = col_double(),
##   SuccDay = col_double(),
##   rndLat = col_double(),
##   alt = col_double()
## )
```

```
data2<-read_csv("data/fric_data_header_2.txt")
```

```
## Parsed with column specification:
## cols(
##   row.index = col_double(),
##   name = col_character(),
##   decimalLongitude = col_double(),
##   decimalLatitude = col_double(),
##   year = col_double(),
##   month = col_double(),
##   day = col_double(),
##   country = col_character(),
##   SuccDay = col_double(),
##   rndLat = col_double(),
##   alt = col_double()
## )
```

```
paste( nrow(data1), "records in format 1;", nrow(data2), "records in format 2")
```

```
## [1] "49243 records in format 1; 233201 records in format 2"
```

```
alldata<-rbind(data1,data2)
rm(data1,data2)
```

##Fric et al includes different species names in results tables than found in data table. In the data curation folder, we match the data names to the results names and create the name_changes.csv file. Here we change names to match results tables:

```
name_changes<-read_csv("data/name_changes.csv")
```

```
## Warning: Missing column names filled in: 'X1' [1]
```

```
## Parsed with column specification:  
## cols(  
##   X1 = col_double(),  
##   result.name = col_character(),  
##   data.name = col_character()  
## )
```

```
table(alldata$name[which(alldata$name %in% name_changes$data_name)])
```

```
## Warning: Unknown or uninitialized column: `data_name`.
```

```
## < table of extent 0 >
```

```
for(namei in 1:nrow(name_changes)) {  
  alldata$name[alldata$name==name_changes$data_name[namei]]<-name_changes$results_name[namei]  
}
```

```
## Warning: Unknown or uninitialized column: `results_name`.
```

```
## Warning: Unknown or uninitialized column: `data_name`.
```

```
## Warning: Unknown or uninitialized column: `results_name`.
```

```
## Warning: Unknown or uninitialized column: `data_name`.
```

```
## Warning: Unknown or uninitialized column: `results_name`.
```

```
## Warning: Unknown or uninitialized column: `data_name`.
```

```
## Warning: Unknown or uninitialized column: `results_name`.
```

```
## Warning: Unknown or uninitialized column: `data_name`.
```

```
## Warning: Unknown or uninitialized column: `results_name`.
```

```
## Warning: Unknown or uninitialized column: `data_name`.
```

```
## Warning: Unknown or uninitialized column: `results_name`.
```

```
## Warning: Unknown or uninitialized column: `data_name`.
```

```
## Warning: Unknown or uninitialized column: `results_name`.
```

```
## Warning: Unknown or uninitialized column: `data_name`.
```

```
## Warning: Unknown or uninitialized column: `results_name`.
```

```
## Warning: Unknown or uninitialized column: `data_name`.
```

```
## Warning: Unknown or uninitialized column: `results_name`.
```

```
## Warning: Unknown or uninitialized column: `data_name`.
```

```
## Warning: Unknown or uninitialized column: `results_name`.
```

```
## Warning: Unknown or uninitialized column: `data_name`.
```

```
## Warning: Unknown or uninitialized column: `results_name`.
```

```
## Warning: Unknown or uninitialized column: `data_name`.
```

```
rm(name_changes)
```

```
## Fric et al identifies datasets by region (N. America, Europe), but the data file does not include this information. We Label data by region using longitude:
```

```
## visualize data density by longitude
```

```
#hist(alldata$decimalLongitude, main="Data density by Longitude")
```

```
#We Label everything East of -40 as Europe, the rest as N. America
```

```
alldata<-alldata %%
```

```
  mutate(region=ifelse(decimalLongitude>=(-40),"Europe","N. America"))
```

```
#Fric et al removed all 1st of month observations and removed one species due to late season nests
```

```
fricdata<-filter(alldata, day!=1, name!="Euphydryas aurinia")
```

```
summary(fricdata)
```

```

##   row.index      name      decimalLongitude      decimalLatitude
## Min.    : 1 Length:275457    Min.   :-162.559    Min.   : 5.787
## 1st Qu.: 2340 Class :character  1st Qu.: -2.676    1st Qu.:52.823
## Median : 7074 Mode  :character  Median : 9.564    Median :55.775
## Mean   :15039                    Mean   : 6.716    Mean   :56.354
## 3rd Qu.:20814                    3rd Qu.: 23.763    3rd Qu.:60.677
## Max.   :85273                    Max.   : 59.333    Max.   :71.216
##
##       year      month      country      day
## Min.    :1616    Min.    : 1.0 Length:275457    Min.    : 2.00
## 1st Qu.:1992    1st Qu.: 6.0 Class  :character  1st Qu.: 9.00
## Median :2002    Median : 7.0 Mode   :character  Median :16.00
## Mean   :1996    Mean   : 6.5                   Mean   :16.19
## 3rd Qu.:2009    3rd Qu.: 7.0                   3rd Qu.:24.00
## Max.   :2015    Max.   :12.0                   Max.   :31.00
## NA's    :57
##       SuccDay      rndLat      alt      region
## Min.    : 2.0    Min.    : 6.00    Min.   :-2666.74 Length:275457
## 1st Qu.:164.0   1st Qu.:53.00   1st Qu.: 23.25 Class  :character
## Median :186.0   Median :56.00   Median : 64.24 Mode   :character
## Mean   :181.2   Mean   :56.29   Mean   : 113.64
## 3rd Qu.:201.0   3rd Qu.:61.00   3rd Qu.: 110.77
## Max.   :361.0   Max.   :71.00   Max.   : 4305.17
##

```

```

#Save formatted and filtered occurrence data used by Fric et al.
save(fricdata,file="data/occurrences_FricAnalysis.RData")

```

*

Data Exploration

Data have now been formatted and filtered to mirror the data used by Fric et al. (2020) and stored into the “occur” tibble.

The following code explores some aspects of the data use in the Fric et al. analysis

```

#
#Here we Visualize data density per regression model from Fric et al.
tempoccurplot<-fricdata %>% group_by(name) %>% tally()
#hist(tempoccurplot$n,xLab="# observations/species", main="Frequency of # observations per species")
#hist(tempoccurplot$n[tempoccurplot$n<5000],xLab="# observations/species up to 5000", xlim=c(0,5000), main="Detail of observations per species", breaks=c(0:5000*100))
rm(tempoccurplot)

#Tally the number of observations per dataset & calculate how each dataset spans latitude, year, altitude
spans.summary<-fricdata %>%
  group_by(name, region) %>%
  add_count(name="fric_n") %>% ## n. records
  group_by(name, region, fric_n) %>%
  summarize(lat_span=(max(rndLat, na.rm=T)-min(rndLat, na.rm=T)),
            year_span=(max(year, na.rm=T)-min(year, na.rm=T)),
            alt_span=round((max(alt, na.rm=T)-min(alt, na.rm=T)),0))

```

```
## `summarise()` regrouping output by 'name', 'region' (override with ` `.groups` argument)
```

```

#calculate # latitudes, onsets, terminations, flight curves = 0
endpt.summary<-fricdata %>%
  group_by(name, region, rndLat) %>%
  # count no. records by latitudinal band
  add_count(name="n_recs") %>%
  #filter to onset & offset dates and label onset dates and offset dates
  filter(SuccDay==min(SuccDay) | SuccDay==max(SuccDay)) %>%
  mutate(onset=ifelse(SuccDay==min(SuccDay),1,0),      term=ifelse(SuccDay==max(SuccDay),1,0)) %>%
  group_by(name, region) %>%
  #create summary statistics by species & region
  summarize(n_lat=length(unique(rndLat)), n_onset=sum(onset), n_term=sum(term), n_flightcurve0s=
sum(n_recs==1) )

```

```
## `summarise()` regrouping output by 'name' (override with ` `.groups` argument)
```

```

#combine summary tables
fric.data.summary<-merge(spans.summary, endpt.summary, by=intersect(names(spans.summary), names(endpt.summary)))
rm(spans.summary)
summary(fric.data.summary)

```

```

##      name          region        fric_n       lat_span
## Length:113     Length:113     Min.   : 4   Min.   : 4.00
## Class :character Class :character 1st Qu.: 71   1st Qu.:23.00
## Mode  :character Mode  :character Median  :184   Median :27.00
##                               Mean   :2438  Mean   :25.93
##                               3rd Qu.:1067 3rd Qu.:30.00
##                               Max.  :51819 Max.  :64.00
##      year_span      alt_span      n_lat       n_onset      n_term
## Min.   : 4.0   Min.   :530   Min.   : 2.00   Min.   : 2.00   Min.   : 2.00
## 1st Qu.:101.0  1st Qu.:2000  1st Qu.:12.00  1st Qu.:14.00  1st Qu.:13.00
## Median :116.0  Median :2653  Median :18.00  Median :19.00  Median :19.00
## Mean   :124.8  Mean   :2678  Mean   :18.61  Mean   :20.09  Mean   :19.81
## 3rd Qu.:138.0  3rd Qu.:3366  3rd Qu.:25.00  3rd Qu.:27.00  3rd Qu.:26.00
## Max.  :399.0   Max.  :5163  Max.  :33.00  Max.  :39.00  Max.  :35.00
##      n_flightcurve0s
## Min.   : 0.000
## 1st Qu.: 2.000
## Median : 3.000
## Mean   : 3.372
## 3rd Qu.: 5.000
## Max.  :10.000

```

*

Explore data by altitude & latitude

This code chunk explores the spatiotemporal representation in the fric.data dataset.

Create Figure 1: Occurrences by altitude & latitude

This code outputs Larsen & Shirey Figure 1, which uses the 4 species presented in Fric et al. Figure 1, to demonstrate the spatiotemporal biases as well as the prevalence of flight periods with a duration of 0 days.

```
summary(fricdata$alt)
```

```

##      Min. 1st Qu. Median      Mean 3rd Qu.      Max.
## -2666.74    23.25   64.24   113.64   110.77  4305.17

```

```

#hist(fricdata$alt)
summary(fricdata$decimalLatitude)

```

```

##      Min. 1st Qu. Median      Mean 3rd Qu.      Max.
##  5.787  52.823  55.775  56.354  60.677  71.216

```

```

#hist(fricdata$decimalLatitude)
##Create Figure 1
#species list
fig1sp<-c("Agriades glandon", "Glaucopsyche lygdamus", "Hesperia comma", "Parnassius smintheus")

#Filter data to these species
fig1data<-fricdata %>%
  filter(name %in% fig1sp)

#Get onset & termination dates (SuccDay)
f1.pheno.data<-fig1data %>%
  group_by(name, region, rndLat) %>%
  mutate(onset=min(SuccDay), term=max(SuccDay), fp=term-onset, singles=ifelse(length(SuccDay)==1, 1, 0))

f1.pheno.data2<-f1.pheno.data %>%
  filter(SuccDay==onset | SuccDay==term)

#A List to store plot panels
tempplot<-list()
fig1panels<-list()

tags<-c("A", "B", "C", "D")

#Create Panels
for(i in 1:2) {
  #paneltitle<-paste(fig1sp[i], "N. America")
  tempplot[[i]] <- ggplot(filter(f1.pheno.data, name==fig1sp[i], region=="N. America"), aes(x=rndLat, y=SuccDay, color=as.factor(singles))) +
    theme_bw() +
    theme(legend.position="none", plot.margin = margin(1,1,1,1, "in")) +
    geom_segment(data=filter(f1.pheno.data2, name==fig1sp[i], region=="N. America"), aes(x=rndLat, y=onset, xend=rndLat, yend=term)) +
    geom_point(aes(color=as.factor(singles))) +
    scale_color_manual(values=c("black", "red")) +
    xlim(min(f1.pheno.data$rndLat), max(f1.pheno.data$rndLat)) + ylim(min(f1.pheno.data$SuccDay), max(f1.pheno.data$SuccDay)) +
    labs(x="Latitudinal Band", y="Day of Year (DOY)", title="") + geom_text(x=min(f1.pheno.data$rndLat), y=max(f1.pheno.data$SuccDay), label=tags[i])

  # with marginal histograms
  fig1panels[[i]] <- ggMarginal(tempplot[[i]], type="histogram")
}

i<-3 #H. comma panel in Fric et al. is from Europe
#paneltitle<-paste(fig1sp[i], "Europe")
tempplot[[i]] <- ggplot(filter(f1.pheno.data, name==fig1sp[i], region=="Europe"), aes(x=rndLat, y=SuccDay, color=as.factor(singles))) +
  theme_bw() +
  theme(legend.position="none", plot.margin = margin(1,1,1,1, "in")) +
  geom_segment(data=filter(f1.pheno.data2, name==fig1sp[i], region=="Europe"), aes(x=rndLat, y=onset, xend=rndLat, yend=term)) +

```

```

geom_point(aes(color=as.factor(singles))) +
scale_color_manual(values=c("black","red")) +
xlim(min(f1.pheno.data$rndLat),max(f1.pheno.data$rndLat)) + ylim(min(f1.pheno.data$SuccDay),
max(f1.pheno.data$SuccDay)) +
labs(x="Latitudinal Band", y="Day of Year (DOY)", title="") + geom_text(x=min(f1.pheno.data
$rndLat), y=max(f1.pheno.data$SuccDay), label=tags[i])

# with marginal histogram
fig1panels[[i]] <- ggMarginal(tempplot[[i]], type="histogram")

##### Figure 1d 2020-07-29 update uses YEAR and DAY to mirror Fric et al.

i<-4
#paneltitle<-paste(fig1sp[i],"N. America")
tempplot[[i]]<- ggplot(filter(f1.pheno.data, name==fig1sp[i], region=="N. America"), aes(x=year,
y=SuccDay, fill=decimalLatitude)) +
geom_point(shape=3) +
theme_bw() +
theme(legend.position="none", plot.margin = margin(1,1,1,1, "in")) +
geom_point(data=filter(f1.pheno.data2, name==fig1sp[i], region=="N. America"), aes(x=year, y=o
nset, fill=decimalLatitude), shape=24) +
geom_point(data=filter(f1.pheno.data, name==fig1sp[i], region=="N. America"), aes(x=year, y=te
rm, fill=decimalLatitude), shape=25) +
scale_fill_gradient(low="azure1", high="black") +
geom_point(data=filter(f1.pheno.data2, name==fig1sp[i], region=="N. America", singles==1), aes
(x=year, y=SuccDay), color="red", shape=16) +
xlim(min(f1.pheno.data$year),max(f1.pheno.data$year)) + ylim(min(f1.pheno.data$SuccDay),max(f
1.pheno.data$SuccDay)) +
labs(x="Year", y="Day of Year (DOY)", title="") + geom_text(x=min(f1.pheno.data$year), y=max(f
1.pheno.data$SuccDay), label=tags[i])

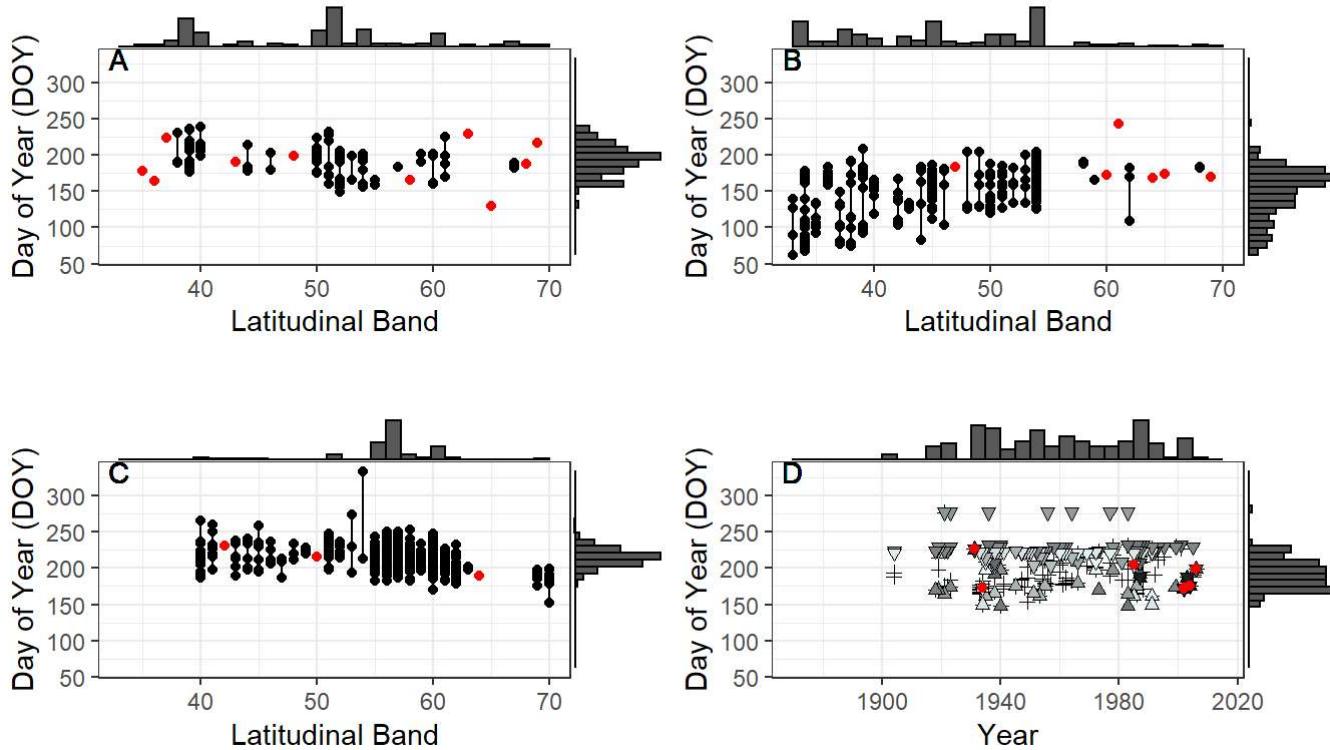
# with marginal histogram
fig1panels[[i]] <- ggMarginal(tempplot[[i]], type="histogram")

grid.arrange(grobs=fig1panels[c(1:4)], nrow=2, ncol=2, top="Visualization of data used in Fric e
t al. for \n (A) Agriades glandon      (B)Glauopsyche lygdamus \n (C)Hesperia comma      (D)Par
nassius smintheus")

```

Visualization of data used in Fric et al. for

- (A) *Agriades glandon* (B) *Glaucopsyche lygdamus*
 (C) *Hesperia comma* (D) *Parnassius smintheus*



```
#Used to create figure 1 pdf
#pdf_filename<-("outputs/Larsen&Shirey2020_FinalFig1.pdf")
#ggsave(pdf_filename, arrangeGrob(grobs=fig1panels[ds1], nrow=2, ncol=2), width=6, height=6, units="in", scale=1,dpi=600)

#rm(f1.pheno.data,f1.pheno.data2, fig1data,fig1panels,fig1sp, ds1,tempplot, tags, i)
```

*

Data curation

Data have now been formatted, identified by region, and summarized.

The following code chunk applies the filters used in the Larsen & Shirey reanalysis and calculates summary data density statistics for all species present in Fric's results to output to Supplemental Table 1.

Our reanalysis excludes datasets along two axes - data density, and voltinism. This code examines data along the data density axis. Unlike Fric et al., we include first day of the month records. We curate raw occurrence data with the following filters prior to estimating phenometrics:

- 1 - remove *Euphydryas aurinia* (as Fric et al. did)
- 2 - altitude in [0m,500m]
- 3 - month in March - November
- 4 - 10 or more records when data is grouped by species, region, year, and latitudinal band

```
#Summarize data availability for Larsen & Shirey re-analysis
#Now, filter data for altitude & for cases with 10 or more records by species-region-year-Latitude
new.data.summary<-alldata %>
  filter(between(alt,0,500), name!="Euphydryas aurinia", month %in% c(3:11)) %>%
  # calculate data availability by species, region, latitude & year
  group_by(name, region, rndLat, year) %>%
  add_count(name="group_n") %>% ## n. observations per group
  filter(group_n>=10) %>% ### filter by 10 or more observations in group
  # calculate reanalysis statistics by species & region
  group_by(name, region) %>%
  add_count(name="curated_n_obs") %>%
  group_by(name, region, curated_n_obs) %>%
  #calculate summary statistics applying data filters
  summarize(curated_n_lat=length(unique(rndLat)), curated_n_fcurve=length(unique(paste(rndLat,year))),,
            curated_lat_span=(max(rndLat, na.rm=T)-min(rndLat, na.rm=T)),
            curated_year_span=(max(year, na.rm=T)-min(year, na.rm=T)),
            curated_alt_span=round((max(alt, na.rm=T)-min(alt, na.rm=T)),0))
```

```
## `summarise()` regrouping output by 'name', 'region' (override with ` `.groups` argument)
```

```
#combine summary tables
supptable1<-merge(fric.data.summary, new.data.summary, by=intersect(names(fric.data.summary), names(new.data.summary)), all.x=T)
head(supptable1)
```

| | | name | region | fric_n | lat_span | year_span | alt_span | n_lat |
|------|--|------------------------|-------------------|------------------|---------------|---------------|------------------|-------|
| ## 1 | | Agriades glandon | N. America | 110 | 34 | 103 | 4042 | 26 |
| ## 2 | | Agriades optilete | Europe | 86 | 25 | 111 | 2940 | 15 |
| ## 3 | | Amblyscirtes vialis | N. America | 88 | 29 | 133 | 2775 | 19 |
| ## 4 | | Anthocharis cardamines | Europe | 31849 | 32 | 168 | 2595 | 33 |
| ## 5 | | Anthocharis sara | N. America | 218 | 28 | 111 | 4417 | 21 |
| ## 6 | | Aphantopus hyperantus | Europe | 30598 | 25 | 399 | 2102 | 26 |
| ## | | n_onset | n_term | n_flightcurve0s | curated_n_obs | curated_n_lat | curated_n_fcurve | |
| ## 1 | | 27 | 27 | 10 | NA | NA | NA | |
| ## 2 | | 15 | 15 | 5 | NA | NA | NA | |
| ## 3 | | 19 | 19 | 7 | NA | NA | NA | |
| ## 4 | | 39 | 35 | 2 | 29134 | 17 | 393 | |
| ## 5 | | 22 | 22 | 6 | NA | NA | NA | |
| ## 6 | | 27 | 28 | 1 | 27879 | 15 | 330 | |
| ## | | curated_lat_span | curated_year_span | curated_alt_span | | | | |
| ## 1 | | NA | | NA | | | | |
| ## 2 | | NA | | NA | | | | |
| ## 3 | | NA | | NA | | | | |
| ## 4 | | 16 | | 80 | | | | |
| ## 5 | | NA | | NA | | | | |
| ## 6 | | 14 | | 79 | | | | |

```
summary(supptable1)
```

```

##      name          region        fric_n       lat_span
## Length:113      Length:113     Min.   : 4   Min.   : 4.00
## Class :character Class :character  1st Qu.: 71   1st Qu.:23.00
## Mode  :character Mode  :character Median  :184   Median :27.00
##                               Mean   :2438   Mean   :25.93
##                               3rd Qu.:1067   3rd Qu.:30.00
##                               Max.  :51819   Max.  :64.00
##
##      year_span      alt_span      n_lat       n_onset      n_term
## Min.   : 4.0   Min.   :530   Min.   : 2.00   Min.   : 2.00   Min.   : 2.00
## 1st Qu.:101.0  1st Qu.:2000  1st Qu.:12.00  1st Qu.:14.00  1st Qu.:13.00
## Median :116.0  Median :2653   Median :18.00   Median :19.00   Median :19.00
## Mean   :124.8   Mean   :2678   Mean   :18.61   Mean   :20.09   Mean   :19.81
## 3rd Qu.:138.0  3rd Qu.:3366  3rd Qu.:25.00  3rd Qu.:27.00  3rd Qu.:26.00
## Max.   :399.0   Max.   :5163   Max.   :33.00   Max.   :39.00   Max.   :35.00
##
##      n_flightcurve0s curated_n_obs curated_n_lat curated_n_fcurve
## Min.   : 0.000   Min.   : 10.0   Min.   : 1.000   Min.   : 1.00
## 1st Qu.: 2.000   1st Qu.: 36.5   1st Qu.: 1.500   1st Qu.: 2.50
## Median : 3.000   Median : 361.0   Median : 3.000   Median :23.00
## Mean   : 3.372   Mean   :3858.6   Mean   : 6.055   Mean   :77.93
## 3rd Qu.: 5.000   3rd Qu.:3928.0   3rd Qu.:10.000   3rd Qu.:124.50
## Max.   :10.000   Max.   :47617.0   Max.   :17.000   Max.   :393.00
## NA's    :58       NA's    :58       NA's    :58       NA's    :58
##
## curated_lat_span curated_year_span curated_alt_span
## Min.   : 0.000   Min.   : 0.00   Min.   : 0.0
## 1st Qu.: 0.500   1st Qu.: 8.00   1st Qu.:214.0
## Median : 5.000   Median : 34.00   Median :379.0
## Mean   : 6.545   Mean   :43.29   Mean   :324.7
## 3rd Qu.:11.500   3rd Qu.:74.00   3rd Qu.:467.5
## Max.   :18.000   Max.   :123.00  Max.   :499.0
## NA's    :58       NA's    :58       NA's    :58

```

```

#output summary table to csv file
#write_csv(supptable1, "Larsen&Shirey_stats_supp_table1.csv")
rm(fric.data.summary, new.data.summary, endpt.summary)

```

*

Data curation 2

This code filters occurrence data for reanalysis by voltinism and data density, and visualizes some differences between datasets curated for the original analysis and this reanalysis. We only include datasets with sufficient data for calculating phenometrics at 3 or more distinct latitudinal bands, so that a linear model can be applied.

```

#FILTER DATA BY VOLTINISM

#get species list without evidence of multiple generations
#Euphydryas aurinia is not included in the voltinism file
voltindata<-read_csv("data/voltinism.csv")

```

```
## Parsed with column specification:  
## cols(  
##   id = col_double(),  
##   name_datafile = col_character(),  
##   name_resultsfile = col_character(),  
##   region = col_character(),  
##   Voltinism = col_character(),  
##   Voltinism_source = col_character(),  
##   `In reanalysis?` = col_double(),  
##   Why_excluded = col_character()  
## )
```

```

voltindata<-na.omit(voltindata[,c(1:8)])
voltindata<-voltindata %>% select(name=name_resultsfile,region,Voltinism)
multi<-c("Bivoltine", "Multivoltine","Sometimes bivoltine","Possible bivoltinism in some subsp."
,"Unconfirmed reports of second brood")
univoltine<-filter(voltindata, !Voltinism %in% multi)
rm(voltindata, multi)

#filter occurrence dataset to these species
reanalysis.data<-merge(alldata, univoltine, by=intersect(names(alldata),names(univoltine)))

#filter data by altitude and data density
reanalysis.data<-reanalysis.data %>%
  filter(between(alt,0,500), month %in% c(3:11)) %>%
  # calculate data availability by species, region, latitude & year
  group_by(name, region, rndLat, year) %>%
  add_count(name="group_n") %>% ## n. observations per group
  filter(group_n>=10) %>% #only groups with at least 10 observations
  group_by(name, region) %>% #group by "dataset"
  mutate(nlat=length(unique(rndLat))) %>% #count how many distinct Latitudinal bands included
  filter(nlat>=3) # need at least 3 Latitudinal bands

#visualize some differences
plotcompar<-list()
plotcompar[[1]]<-ggplot(data=fricdata, aes(x=region, y=alt) ) +
  geom_boxplot(outlier.colour="red", outlier.shape=16, outlier.size=2, notch=FALSE) + ggtitle(label="Original dataset altitudes")

plotcompar[[2]]<-ggplot(data=reanalysis.data, aes(x=region, y=alt) ) +
  geom_boxplot(outlier.colour="red", outlier.shape=16, outlier.size=2, notch=FALSE) + ggtitle(label="Reanalysis dataset altitudes") + ylim(min(fricdata$alt),max(fricdata$alt))

plotcompar[[3]]<-ggplot(data=fricdata, aes(x=region, y=rndLat) ) +
  geom_boxplot(outlier.colour="red", outlier.shape=16, outlier.size=2, notch=FALSE) + ggtitle(label="Original dataset latitudes")

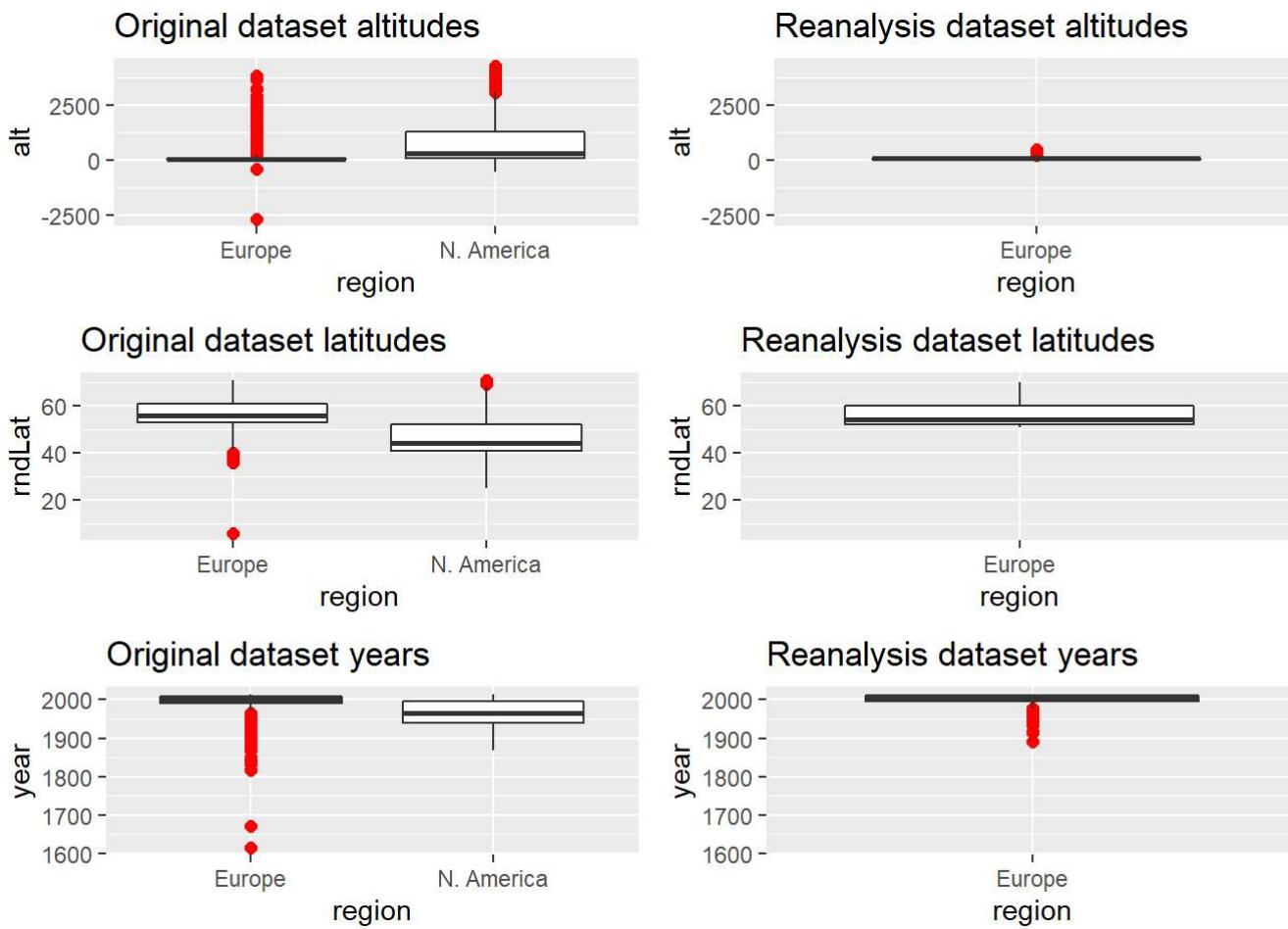
plotcompar[[4]]<-ggplot(data=reanalysis.data, aes(x=region, y=rndLat) ) +
  geom_boxplot(outlier.colour="red", outlier.shape=16, outlier.size=2, notch=FALSE) + ggtitle(label="Reanalysis dataset latitudes") + ylim(min(fricdata$rndLat), max(fricdata$rndLat))

plotcompar[[5]]<-ggplot(data=filter(fricdata, !is.na(year)), aes(x=region, y=year) ) +
  geom_boxplot(outlier.colour="red", outlier.shape=16, outlier.size=2, notch=FALSE) + ggtitle(label="Original dataset years")

plotcompar[[6]]<-ggplot(data=reanalysis.data, aes(x=region, y=year) ) +
  geom_boxplot(outlier.colour="red", outlier.shape=16, outlier.size=2, notch=FALSE) + ggtitle(label="Reanalysis dataset years") + ylim(min(fricdata$year, na.rm=T), max(fricdata$year, na.rm=T))

grid.arrange(grobs=plotcompar[c(1:6)], nrow=3)

```



*

Estimate phenometrics using phest

This chunk of code estimates onset and offset phenometrics by species-region-year-latitudinal_band using curated data.

We use the phest package to estimate onset and offset of flight periods based on occurrence data, when at least 10 observations exist for a species-region-year-latitudinal_band unit. The phest package applies a weibull distribution. Please note that this chunk does take a few minutes to run. Also, warnings are automatically generated by “phest” when a correction is applied to the phenometric estimate. Additionally, “phest” throws a warning for CI estimation. We have explored these warnings and don’t believe that there is any problem continuing with the estimates produced; therefore we have suppressed the warning messages here.

```

rm(plotcompar)

#For each species & region, calculate phenometrics
datasets<-reanalysis.data %>% group_by(name, region) %>% tally()
pheno.est<-data.frame(name=character(0),region=character(0),year=integer(0),rndLat=integer(0),on
set.est=numeric(0),onset.low=numeric(0),onset.high=numeric(0),offset.est=numeric(0),offset.low=n
umeric(0),offset.high=numeric(0))

for(rowi in 1:nrow(datasets)){ # for each unique dataset
  namei<-datasets$name[rowi]
  regi<-datasets$region[rowi]
  index <- 1 # create/reset an indexer
  pheno.estimates <- list() # create/refresh a blank list per group
  rowi.data<-filter(reanalysis.data, name==namei, region==regi)
  for(yr in unique(rowi.data$year)){ # and each unique year
    for(lat in unique(rowi.data$rndLat)){ # and each unique latitude
      temp <- filter(rowi.data, rndLat==lat, year==yr) # filter the occurrence data for each gro
up

      if(nrow(temp) > 9){ # if there are at least 10 occurrences, then...
        estimates <- c(namei, regi, yr, lat, nrow(temp),
                         suppressWarnings(weib.limit(temp$SuccDay, upper=FALSE, alpha=0.05)),
                         suppressWarnings(weib.limit(temp$SuccDay, upper=TRUE, alpha=0.05))) # cal
culate estimates for the group: onset, offset
        pheno.estimates[[index]] <- estimates # shuttle those into a list
        index <- index+1
      } #end if enough occurrences
    } #end Lat
  } #end yr
  df <- data.frame(matrix(unlist(pheno.estimates), nrow=length(pheno.estimates), byrow=TRUE),str
ingsAsFactors=FALSE)
  names(df)<-c("name","region","year","rndLat","n","onset.est","onset.low","onset.high","offset.
est","offset.low","offset.high")
  pheno.est<-rbind(pheno.est, df)
}
for(col in 3:11) {
  pheno.est[,col]<-as.numeric(pheno.est[,col])
}

#Format & store data
pheno.data<-pheno.est %>%
  mutate(unit=paste(name, rndLat, year,sep="-")) %>%
  select(unit,onset.est,offset.est,name,region,rndLat,year,n) %>%
  mutate(onset=round(onset.est,0),term=round(offset.est,0))
pheno.data<-na.omit(pheno.data)
#Weibull estimator doesn't bound so
#We bounded all onset & termination metrics y [60,330], Limiting flight periods to March - Novem
ber
pheno.data$onset[pheno.data$onset<60]<-60
pheno.data$term[pheno.data$term>330]<-330

#save(pheno.data, file="data/phenometrics.RData")

```

*

Statistical models for phenometrics

This code uses estimated onset and offset phenometrics in linear models to examine phenological patterns with latitude and year. Other statistical models may be more appropriate for a de novo analysis, but here we want our statistical model to parallel the Fric et al. model in intention, but using multiple regression instead of residual regression.

```

#If we want to skip phest and phenometric estimation:
#Load("data/phenometrics.RData")

datasets<-pheno.data %>%
  group_by(name, region) %>%
  tally()
pheno.data<-na.omit(pheno.data)
#Loop through datasets, run model for phenology by species & region, and store LM parameters
onsetpheno<-list()
termpheno<-list()
onset1<-NULL
term1<-NULL
axes<-NULL

for(rowi in 1:nrow(datasets)) {
  pheno.rowi<-pheno.data %>%
    filter(name==datasets$name[rowi], region==datasets$region[rowi])
#estimate model params for onset
  onset.lm<-summary(lm(onset~rndLat+year, data=pheno.rowi))$coefficients #estimate model params for termination
  term.lm<-summary(lm(term~rndLat+year, data=pheno.rowi))$coefficients
#store
  onsetpheno[[rowi]]<-onset.lm
  termpheno[[rowi]]<-term.lm

  #onset
  temponset<-matrix(unlist(onset.lm[c(2:3),]), ncol=4, byrow=F)
  onset1<-rbind(onset1, temponset)
  axes<-c(axes, row.names(onset.lm)[c(2:3)])
#termination
  tempterm<-matrix(unlist(term.lm[c(2:3),]), ncol=4, byrow=F)
  term1<-rbind(term1, tempterm)
  rm(pheno.rowi, onset.lm, term.lm, temponset, tempterm)
}

#Create results dataframes: onset
onset1<-as.data.frame(onset1)
colnames(onset1)<-c("param.est", "param.se", "param.t", "param.p")
onset1$param<-axes
onset1$metric<-"onset"
onset1$name<-rep(datasets$name, each=2)
onset1$region<-rep(datasets$region, each=2)
onset1$n<-rep(datasets$n, each=2)

#Create results dataframes: termination
term1<-as.data.frame(term1)
colnames(term1)<-c("param.est", "param.se", "param.t", "param.p")
term1$param<-axes
term1$metric<-"termination"
term1$name<-rep(datasets$name, each=2)
term1$region<-rep(datasets$region, each=2)
term1$n<-rep(datasets$n, each=2)

```

```

result<-as.data.frame(rbind(onset1, term1))
result<-result %>%
  mutate(response=ifelse(param.p<0.05,ifelse(param.est>0,1,-1),0))

#Plot coefficients colored by response sign
#coef.plot1<-ggplot(data=filter(result, param=="rndLat"), aes(x=as.factor(metric), y=param.est))
+
# geom_boxplot() +
# geom_jitter(data=filter(result, param=="rndLat"), aes(x=as.factor(metric), y=param.est, color=as.factor(response)), width=0.2, height=0, shape=17) +
# labs(x="", y="Latitude coefficient") +
# scale_color_manual(values=c("blue", "darkgray", "darkgreen")) +
# theme_light() + theme(Legend.position = "none")
#coef.plot1
#The positive and neutral slopes of onset ~ latitude indicate that species typically emerge either at similar or later times of year at higher latitudes. This along with the varied slopes of termination ~ latitude indicate that most flight periods may generally shift later and/or have shorter duration at higher latitudes.

#NOT in manuscript but exploratory: Using only coefficients and significance without confidence intervals, what phenological patterns are present?
slopediff<-NULL
for(spi in unique(result$name)) {
  d.start<-ifelse(filter(result,param=="rndLat",name==spi,metric=="onset")$response>0,"later","same")
  d.duration<- ifelse(filter(term1,param=="rndLat",name==spi)$param.est-filter(onset1,param=="rndLat",name==spi)$param.est<0,"shorter",ifelse(filter(term1,param=="rndLat",name==spi)$param.est-filter(onset1,param=="rndLat",name==spi)$param.est>0,"longer","same"))
  slopediff<-c(slopediff, paste(d.start,d.duration,sep="."))
}
table(slopediff)

```

```

## slopediff
##   later.longer later.shorter same.longer same.shorter
##                 1              12               2               5

```

*

Compare statistical results to Fric et al.

This code uses model outputs and compares them to the results of the Fric et al. analysis. It outputs Figure 2.

```

##Results and visualizations

fric.results.lat<-read_excel("data/Fric_results.xlsx", sheet="models")
na.omit(read_excel("fric_supplements/Supplementary Table 2_final.xlsx", sheet="~latitude", range ="A3:A113"))

```

```
## # A tibble: 105 x 1
##   Species
##   <chr>
## 1 Carterocephalus palaemon (Pallas, 1771)
## 2 Hesperia comma (Linnaeus, 1758)
## 3 Thorybes pylades (Scudder, 1870)
## 4 Erynnis icelus (Scudder & Burgess, 1870)
## 5 Erynnis persius (Scudder, 1863)
## 6 Pyrgus centaureae (Rambur, 1840)
## 7 Amblyscirtes vialis (W.H. Edwards, 1862)
## 8 Parnassius smintheus Doubleday, 1847
## 9 Papilio machaon Linnaeus, 1758
## 10 Papilio canadensis Rothschild & Jordan, 1906
## # ... with 95 more rows
```

```

datasets$set<-paste(datasets$name,datasets$region,sep="-")
fric.results<-fric.results.lat %>%
  mutate(region=ifelse(region=="NA","N. America","Europe")) %>%
  mutate(reanalyzed=ifelse(paste(name,region,sep="-") %in% datasets$set,1,0))
#Model 1 = Fric Direct regression, all species
fric1<-fric.results %>%
  filter(model=="~latitude") %>%
  select(name,region,onset_coef, onset_response, term_coef, term_response) %>%
  mutate(modelnum=1)

#Model 3 = Fric Direct regression, reanalyzed species
fric3<-fric.results %>%
  filter(model=="~latitude", reanalyzed==1) %>%
  select(name,region,onset_coef, onset_response, term_coef, term_response) %>%
  mutate(modelnum=3)

#Model 2 = Fric residual regression, all species
fric2<-fric.results %>%
  filter(model=="~latitude|altitude+year") %>%
  select(name,region,onset_coef, onset_response, term_coef, term_response) %>%
  mutate(modelnum=2)

#Model 4 = Fric residual regression, reanalyzed species
fric4<-fric.results %>%
  filter(model=="~latitude|altitude+year", reanalyzed==1) %>%
  select(name,region,onset_coef, onset_response, term_coef, term_response) %>%
  mutate(modelnum=4)

#Model 5 = Reanalysis multiple regression
temp<-pivot_wider(filter(result, param=="rndLat"), id_cols =c(name, region),names_from=metric,values_from=c(param.est,param.p, response) )

result5<-temp %>%
  select(name, region, onset_coef=param.est_onset, onset_response=response_onset,term_coef=param.est_termination,term_response=response_termination) %>%
  mutate(modelnum=5)
rm(temp,fric.results.lat)

#Combine all results into 1 data frame
result.compar<-as.data.frame(rbind(fric1,fric2,fric3,fric4,result5))
result.compar$modelnum<-as.factor(result.compar$modelnum)
result.compar$s1<-1
##Create Figure 2
colorscheme<-c("blue", "darkgray", "darkgreen")
modelnames<-c("AllFric1","AllFric2","Fric1","Fric2","Reanalysis")
#Panels A, D: compare coefficients
#Panel A: Onset coefficients
onset.sp<-ggplot(data=filter(result.compar, as.numeric(modelnum)>3), aes(x=name, y=onset_coef, shape=as.factor(modelnum), fill=as.factor(onset_response))) +
  geom_point(color="black") +
  scale_shape_manual(values=c(22,21)) +
  scale_fill_manual(values=c("white","black")) +
  geom_hline(yintercept=0) +

```

```

scale_y_continuous(breaks=seq(-8,8,2)) +
  labs(x="", y="Latitude coefficient") + coord_flip() +
  theme_light() + theme(legend.position = "none", axis.text=element_text(size=8))
#onset.sp
#Panel D: Termination coefficients
term.sp<-ggplot(data=filter(result.compar, as.numeric(modelnum)>3), aes(x=name, y=term_coef, sha
pe=as.factor(modelnum), fill=as.factor(term_response))) +
  geom_point(color="black") +
  scale_shape_manual(values=c(22,21)) +
  scale_fill_manual(values=c("black","white","black")) +
  geom_hline(yintercept=0) +
  scale_y_continuous(breaks=seq(-8,8,2)) +
  labs(x="", y="Latitude coefficient") + coord_flip() +
  theme_light() + theme(legend.position = "none", axis.text=element_text(size=8))
#term.sp

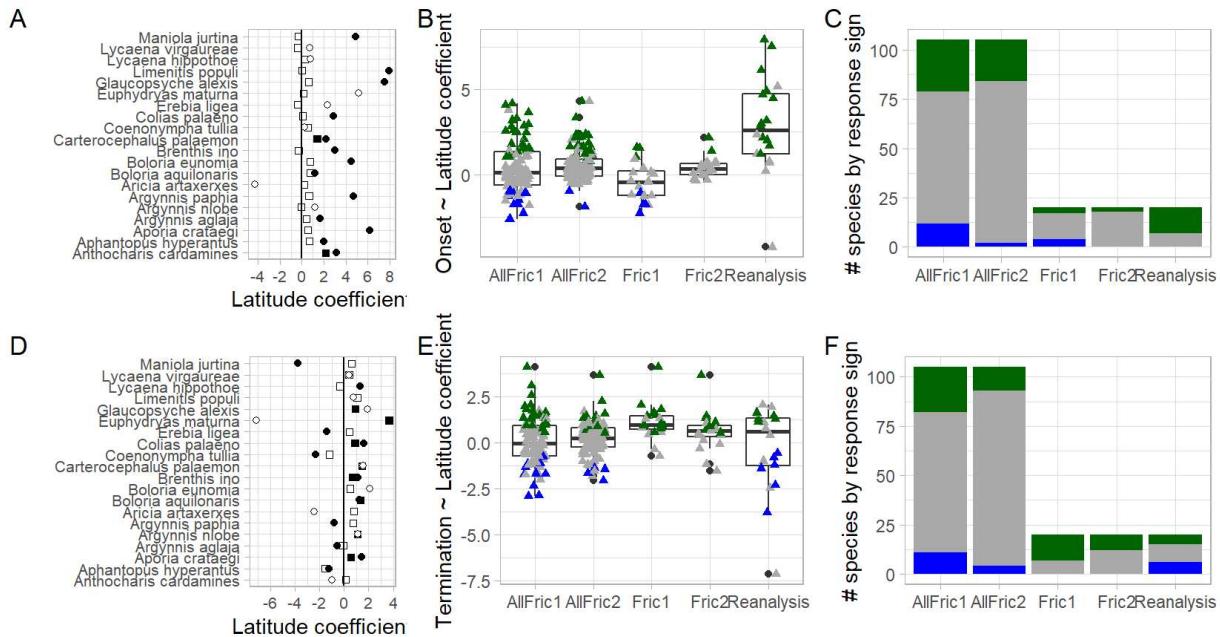
#Panels B, E: response boxplots
#Panel B: Onset
onset.c<-ggplot(data=result.compar, aes(x=modelnum, y=onset_coef)) +
  geom_boxplot() +
  geom_jitter(data=filter(result.compar), aes(x=modelnum, y=onset_coef, color=as.factor(onset_r
esponse)), width=0.2, height=0, shape=17) +
  labs(x="", y="Onset ~ Latitude coefficient") +
  scale_color_manual(values=colorscheme) +
  theme_light() + theme(legend.position = "none") +
  scale_x_discrete(breaks=c(1:5),labels=modelnames)
#onset.c
#Panel D: termination
term.c<-ggplot(data=result.compar, aes(x=modelnum, y=term_coef)) +
  geom_boxplot() +
  geom_jitter(data=filter(result.compar), aes(x=modelnum, y=term_coef, color=as.factor(term_res
ponse)), width=0.2, height=0, shape=17) +
  labs(x="", y="Termination ~ Latitude coefficient") +
  scale_color_manual(values=colorscheme) +
  theme_light() + theme(legend.position = "none") +
  scale_x_discrete(breaks=c(1:5),labels=modelnames)

#Panels C, F: stacked barplots
#Panel c: Onset responses
onset.st<-ggplot(data=result.compar, aes(x=modelnum, y=s1, fill=as.factor(onset_response))) +
  geom_bar(position=position_stack(reverse=T), stat="identity") +
  scale_fill_manual(values=colorscheme) +
  labs(x="", y="# species by response sign") + theme_light() + theme(legend.position = "none") +
  scale_x_discrete(breaks=c(1:5),labels=modelnames)
#Panel F: Termination responses
term.st<-ggplot(data=result.compar, aes(x=modelnum, y=s1, fill=as.factor(term_response))) +
  geom_bar(position=position_stack(reverse=T), stat="identity") +
  scale_fill_manual(values=colorscheme) +
  theme_light() +
  labs(x="", y="# species by response sign") + theme(legend.position = "none") +
  scale_x_discrete(breaks=c(1:5),labels=modelnames)
#term.st

```

```
##Combine panels into Figure 2:
p1<-onset.sp+labs(tag="A")
p2<-onset.c+labs(tag="B")
p3<-onset.st+labs(tag="C")
p4<-term.sp+labs(tag="D")
p5<-term.c+labs(tag="E")
p6<-term.st+labs(tag="F")

#pdf_filename<-("output/LarsenShirey2020_Fig2.pdf")
fig2<-grid.arrange(ncol=3, grobs=list(p1, p2, p3, p4, p5, p6), top="\n\n", bottom="\n\n", left="\n\n", right="\n\n", width=10, height=5)
```



```
#ggsave(pdf_filename, arrangeGrob(fig2, nrow=1), width=10, height=5, scale=1.5, dpi=600, units="in")  
fig2
```

```
## TableGrob (4 x 5) "arrange": 10 grobs
##   z   cells   name      grob
## 1 1 (2-2,2-2) arrange    gtable[layout]
## 2 2 (2-2,3-3) arrange    gtable[layout]
## 3 3 (2-2,4-4) arrange    gtable[layout]
## 4 4 (3-3,2-2) arrange    gtable[layout]
## 5 5 (3-3,3-3) arrange    gtable[layout]
## 6 6 (3-3,4-4) arrange    gtable[layout]
## 7 7 (1-1,2-4) arrange text[GRID.text.1146]
## 8 8 (4-4,2-4) arrange text[GRID.text.1147]
## 9 9 (1-4,1-1) arrange text[GRID.text.1148]
## 10 10 (1-4,5-5) arrange text[GRID.text.1149]
```

*

Create statistics for results table (Supplemental Table 2)

This code outputs a results table that is a partial Supplemental Table 2 - it is currently missing the 'year' analyses from Fric et al., as our focus is on the latitudinal patterns.

```

fric.table<-fric.results %>%
  filter(reanalyzed==1) %>%
  select(name, region, model, onset_response, onset_p, onset_coef, term_response, term_p=term_p_mean, term_coef)

fric.onset<-fric.table %>%
  select(name:model, Fric_singleRegression_sign=onset_response, Fric_singleRegression_p=onset_p,
Fric_singleRegression_coef=onset_coef,) %>%
  mutate(phenometric="onset")
fric.term<-fric.table %>%
  select(name:model, Fric_singleRegression_sign=term_response, Fric_singleRegression_p=term_p,Fric_singleRegression_coef=term_coef,) %>%
  mutate(phenometric="termination")
fric.table<-rbind(fric.onset, fric.term)
fric.table<-fric.table %>%
  mutate(indep.variable=ifelse(model=="~latitude","latitude",ifelse(model=="~year","year",0)) )

fric1<-filter(fric.table,indep.variable %in% c("latitude","year"))
fric2<-fric.table %>%
  filter(!indep.variable %in% c("latitude","year")) %>%
  select(name, region,phenometricFric_resid.regress_sign=Fric_singleRegression_sign, Fric_resid.regress_p=Fric_singleRegression_p,Fric_resid.regress_coef=Fric_singleRegression_coef) %>%
  mutate(indep.variable="latitude")

reanalysis.table<-result %>%
  select(name, region, phenometric=metric, param, response, param.p, param.est) %>%
  mutate(indep.variable=ifelse(param=="year","year","latitude"))%>%
  select(name, region, phenometric, indep.variable, Reanalysis_sign=response, Reanalysis_p=param.p, Reanalysis_coef=param.est)

supptable2<-merge(reanalysis.table,fric1[,c(1,2,4:8)], by=intersect(names(reanalysis.table),names(fric1[,c(1,2,4:8)])))

supptable2<-merge(supptable2,fric2,by=intersect(names(supptable2),names(fric2))), all.x=T)
summary(supptable2)

```

```

##      name          region      indep.variable      phenometric
## Length:120      Length:120      Length:120      Length:120
## Class :character Class :character Class :character Class :character
## Mode  :character Mode  :character Mode  :character Mode  :character
##
##
##
##
##  Reanalysis_sign  Reanalysis_p   Reanalysis_coef
## Min.  :-1.0000    Min.  :0.0000000  Min.  :-7.1347
## 1st Qu.: 0.0000   1st Qu.:0.003991  1st Qu.:-0.1605
## Median : 0.0000   Median :0.048235  Median : 0.4745
## Mean   : 0.1833   Mean   :0.190249  Mean   : 0.8932
## 3rd Qu.: 1.0000   3rd Qu.:0.265633  3rd Qu.: 1.9223
## Max.   : 1.0000   Max.   :0.940501  Max.   : 7.9281
##
## Fric_singleRegression_sign Fric_singleRegression_p Fric_singleRegression_coef
## Min.  :-1.0000        Min.  :0.0000006  Min.  :-2.2694
## 1st Qu.: 0.0000        1st Qu.:0.0161471  1st Qu.:-0.3890
## Median : 0.0000        Median :0.1161413  Median : 0.1626
## Mean   : 0.2167        Mean   :0.2768246  Mean   : 0.2420
## 3rd Qu.: 1.0000        3rd Qu.:0.5125612  3rd Qu.: 0.9424
## Max.   : 1.0000        Max.   :0.9870071  Max.   : 4.1062
##
## phenometricFric_resid.regress_sign Fric_resid.regress_p
## Min.  :0.00          Min.  :0.00013
## 1st Qu.:0.00          1st Qu.:0.05282
## Median :0.00          Median :0.27878
## Mean   :0.25          Mean   :0.36836
## 3rd Qu.:0.25          3rd Qu.:0.64797
## Max.   :1.00          Max.   :0.99961
## NA's   :40            NA's   :40
##
## Fric_resid.regress_coef
## Min.  :-1.52542
## 1st Qu.: 0.06082
## Median : 0.53580
## Mean   : 0.49531
## 3rd Qu.: 0.75339
## Max.   : 3.66029
## NA's   :40

```

##This partial supplementary table 2 does not include the residual regression year results, which would fill in the NA's in the table.

```
#write.csv(supptable2,file="output/supp_table2_part.csv")
```

*

Create panels for Supplemental Figure 1

In this code chunk, we previously used `lm.model$call` references in `geom_smooth`, which created a string of outputs showing the calls. The current simple `lm` still includes `geom_smooth` output text, which would be nice to suppress if we can figure that out.

*

A small break between creating the panels and assmembling Supplemental Figure 1. The chunk below combines the panels into Supplemental Figure 1.

*

Below is the code used to create Supplemental Figure 1 in R, for documentation.

We hope to add a live link to View Supplemental Figure 1 pdf.

This is the end of this analysis.

Author notes - Future updates should: Remove variables when we're done with them See if we can suppress `geom_smooth()` messages