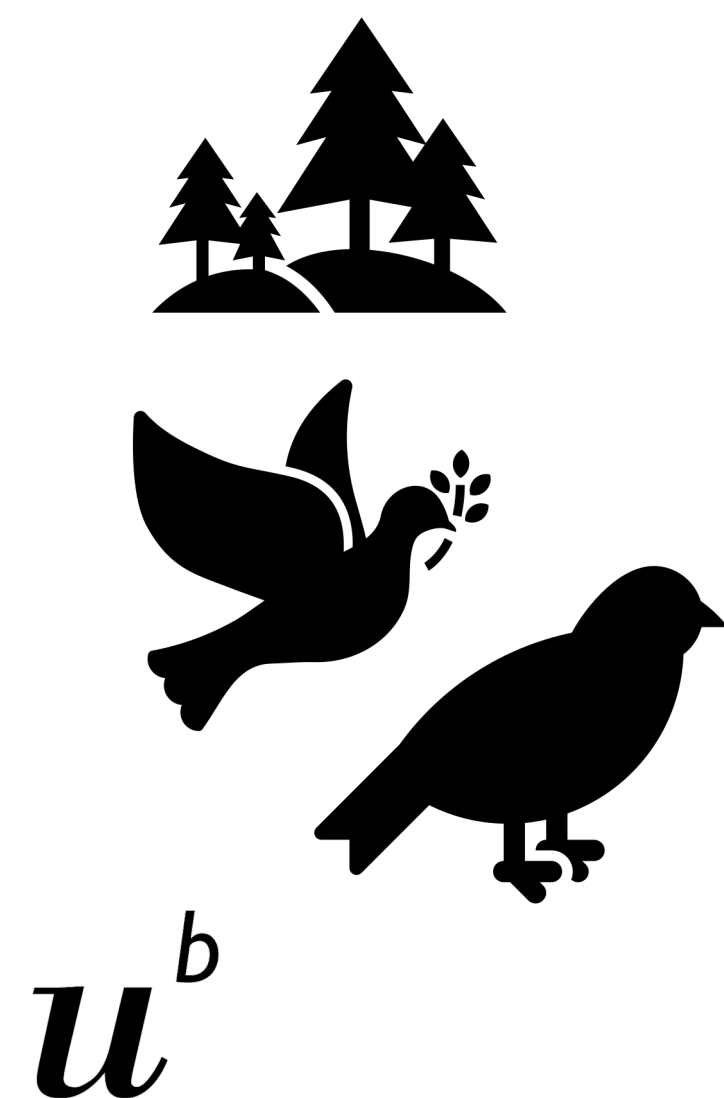


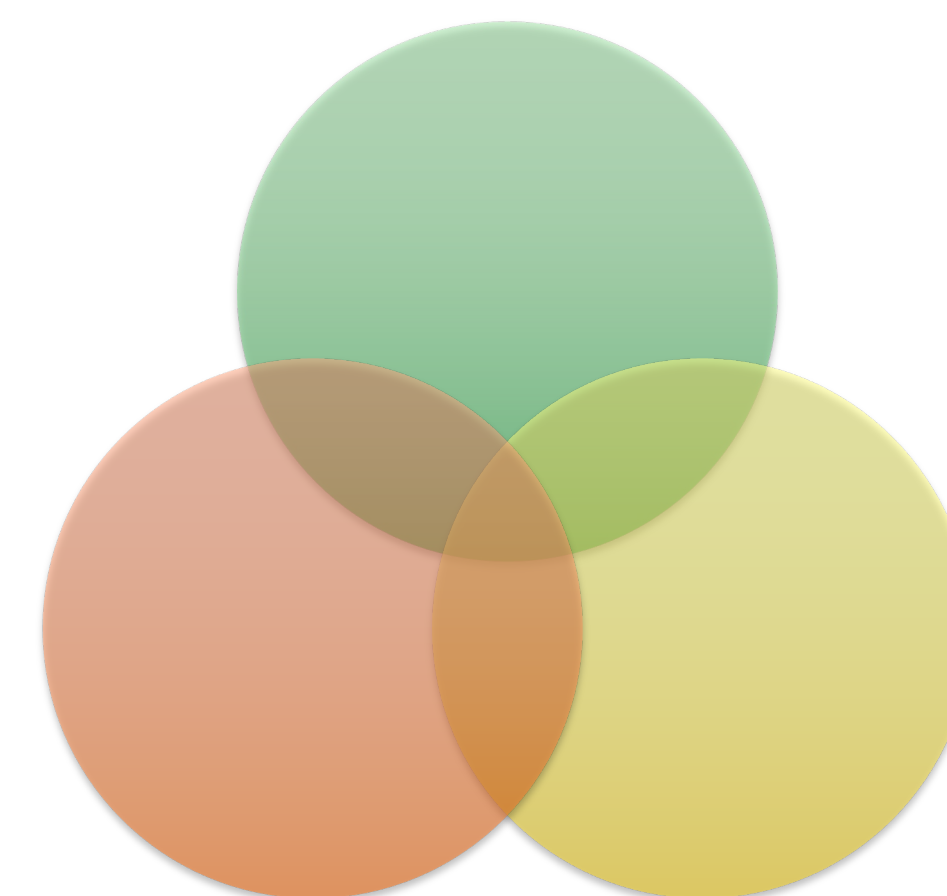
EXTREME-VALUE MODELLING OF MIGRATORY BIRD ARRIVAL DATES: INSIGHTS FROM CITIZEN SCIENCE DATA

Jonathan Koh, Thomas Opitz



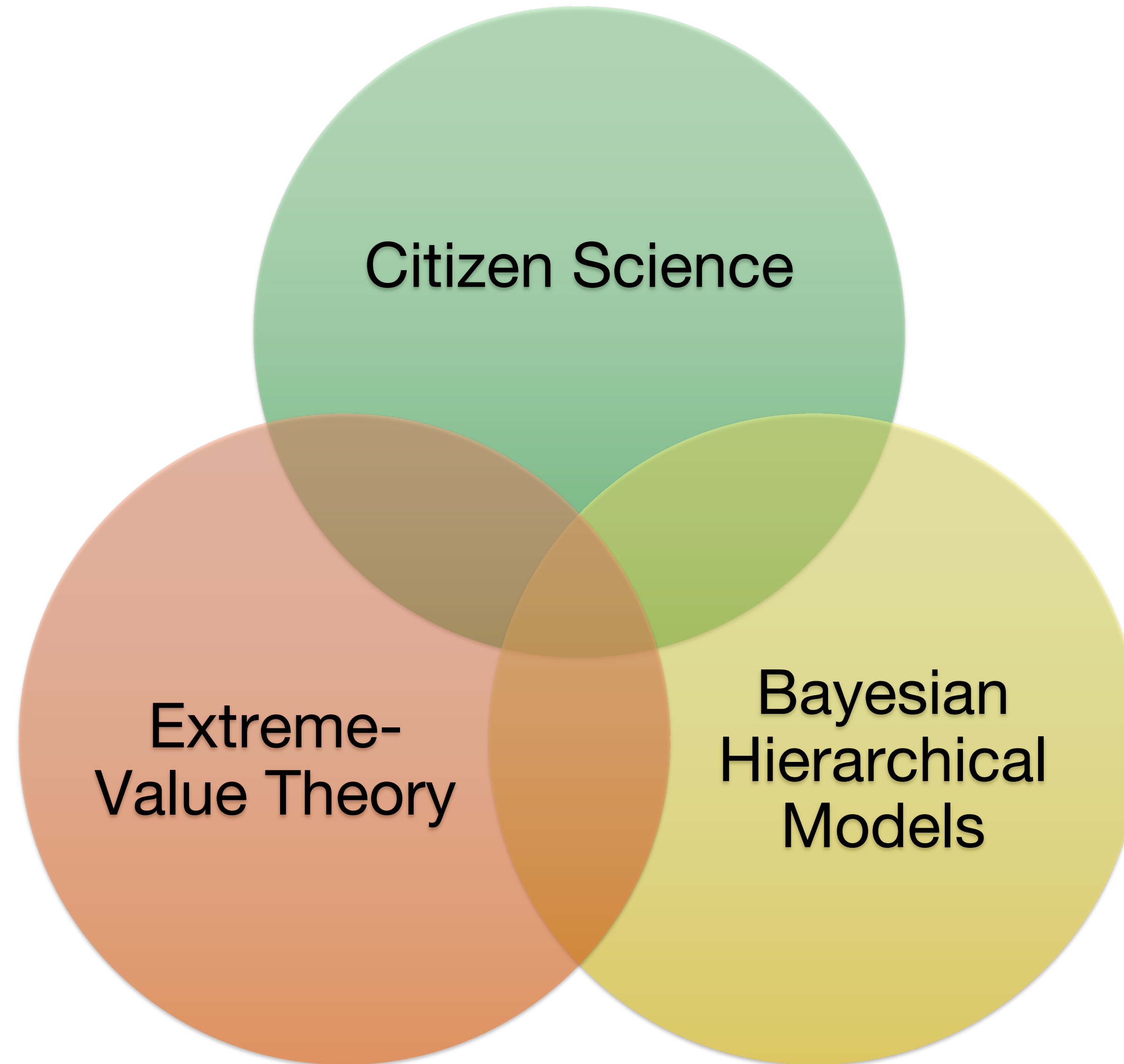
^b
UNIVERSITÄT
BERN

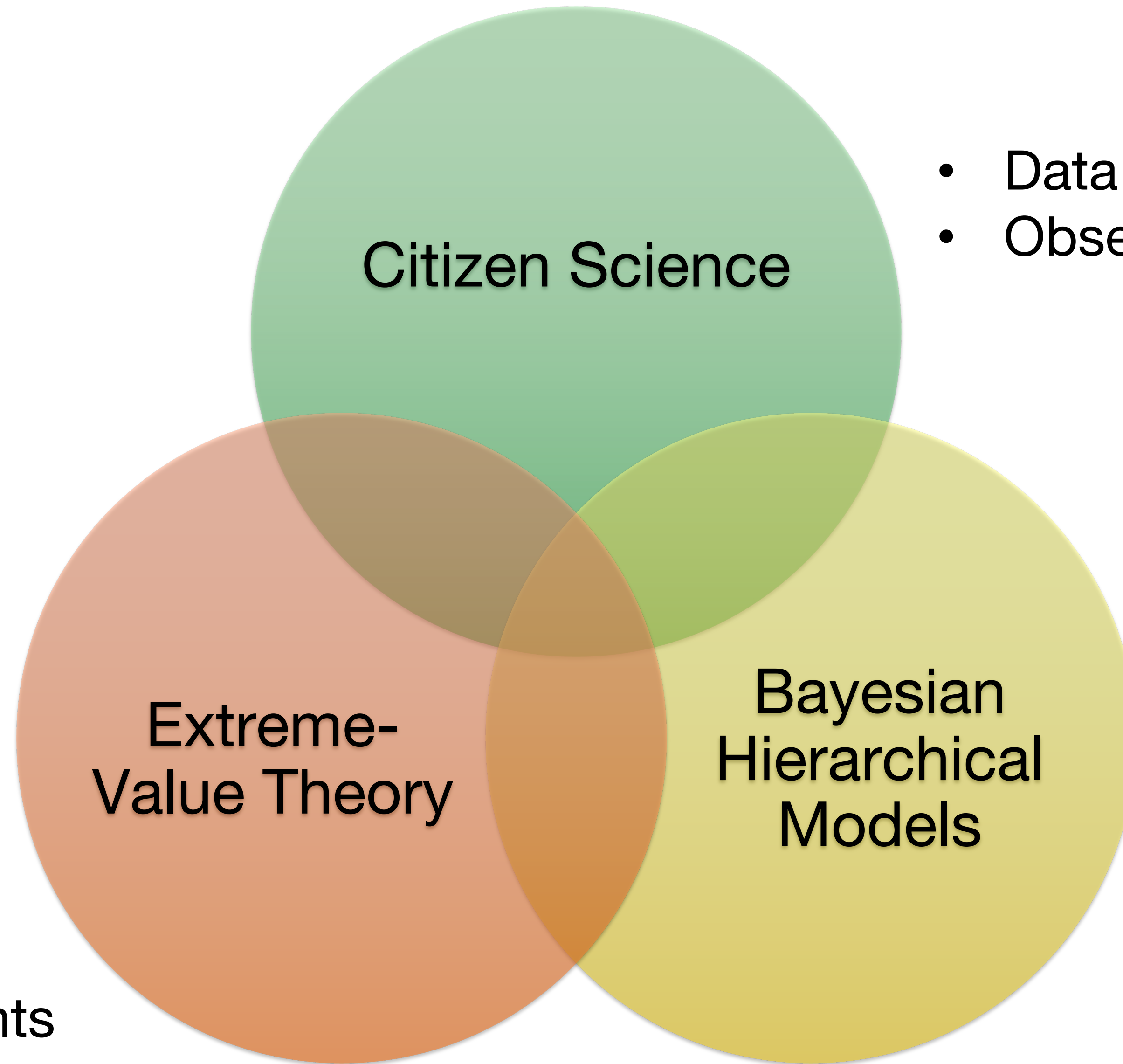
OESCHGER CENTRE
CLIMATE CHANGE RESEARCH



INRAE

RSS meeting 2024
Discussion paper session, 03/09/2024





- Data fusion
- Observational bias

- Extremes of phenological events

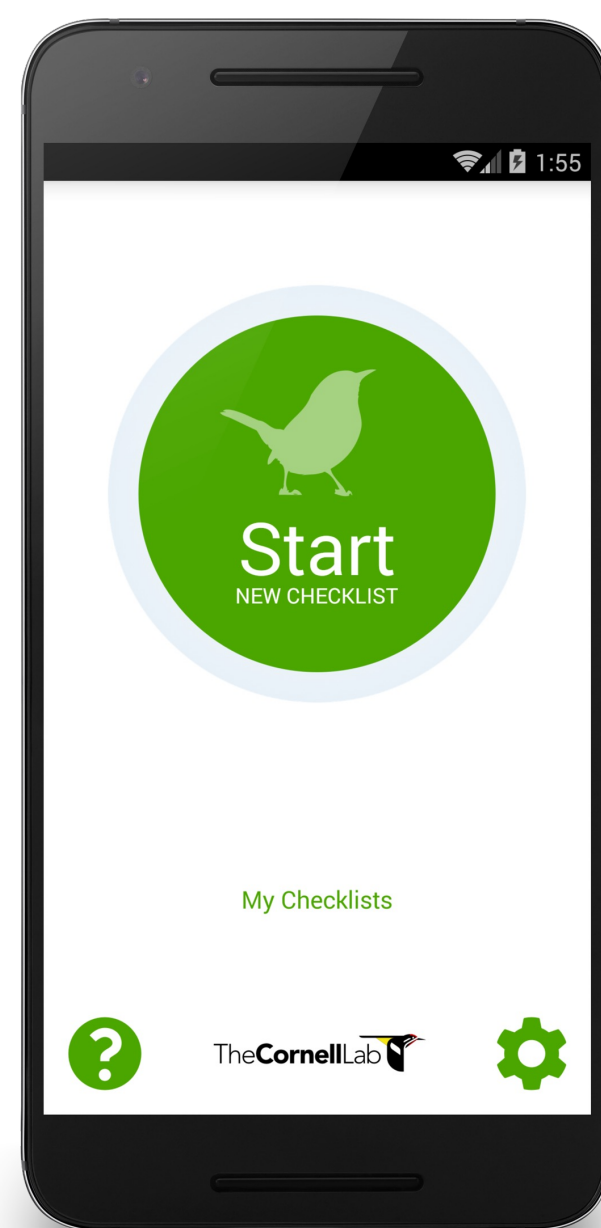
- Inference on important latent processes

How does eBird work?

1. Dr. Andrew Garrett is a birder
2. While hiking, he opens the eBird app and starts a `checklist`. The app notes the date and time he starts birding, where he has travelled during the checklist and how long he has been birding

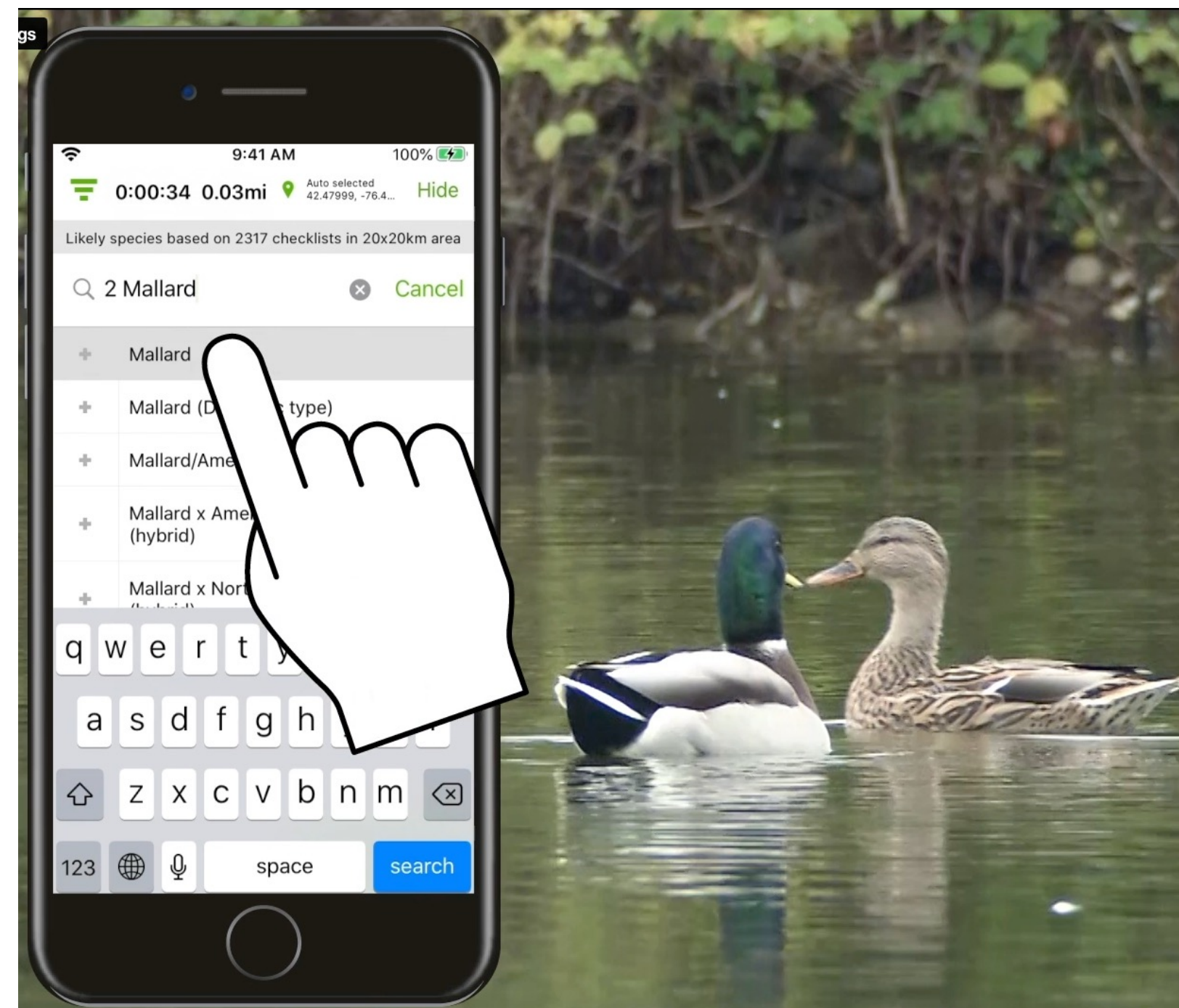
How does eBird work?

1. Dr. Andrew Garrett is a birder
2. While hiking, he opens the eBird app and starts a `checklist`. The app notes the date and time he starts birding, where he has travelled during the checklist and how long he has been birding



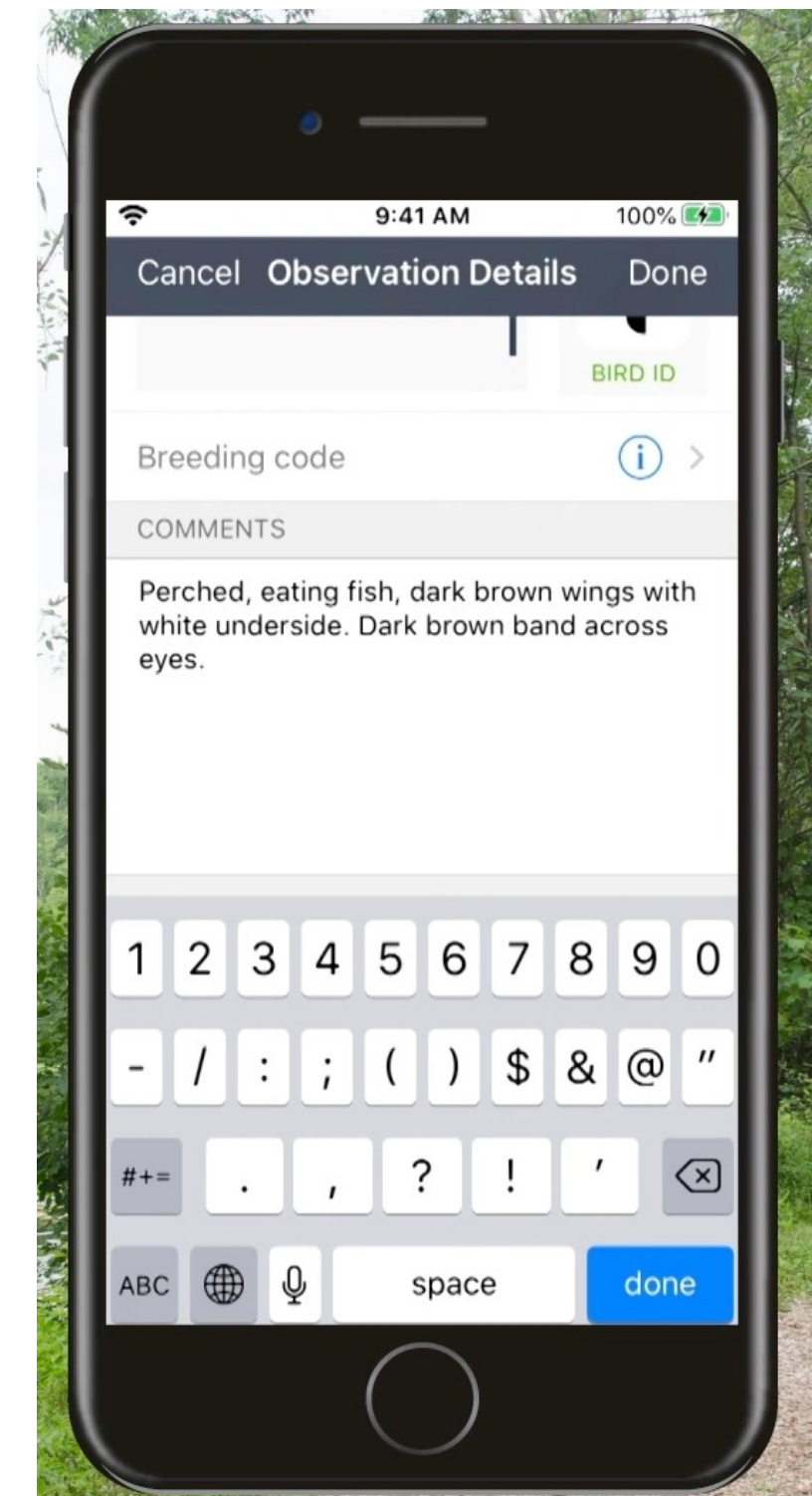
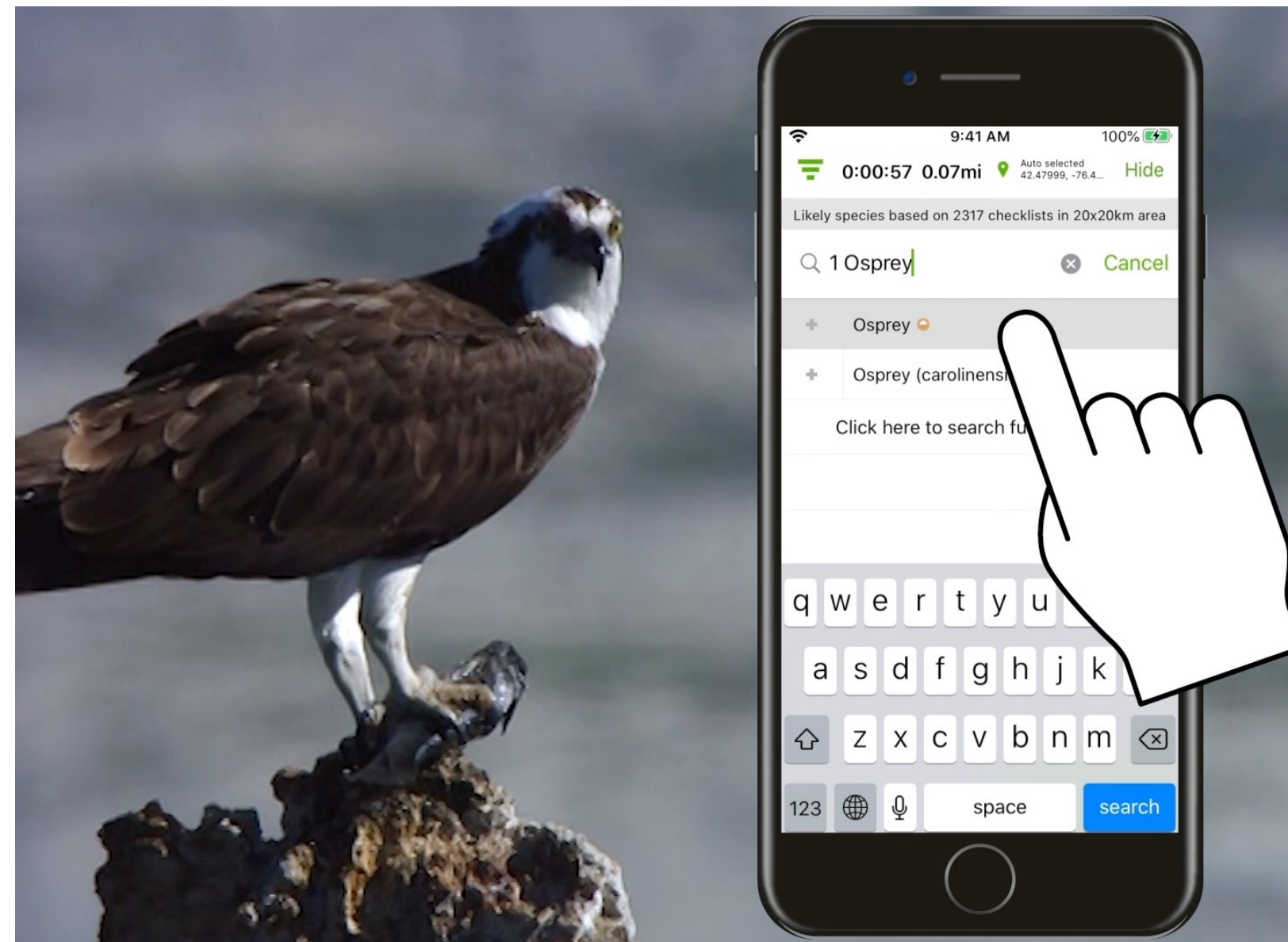
How does eBird work?

3. He spots a Mallard, and can easily record it in the app (based on a pack of recommended bird species)



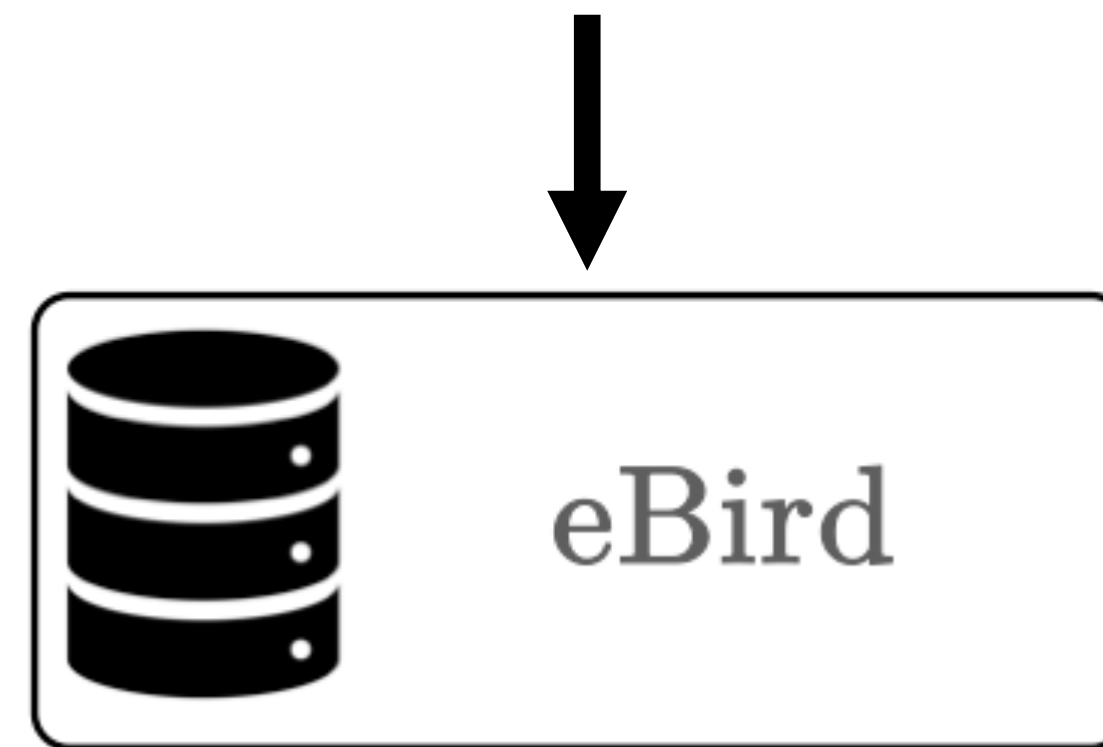
How does eBird work?

4. He spots an Osprey, and records it in the app. He can also supply more information (including media files)



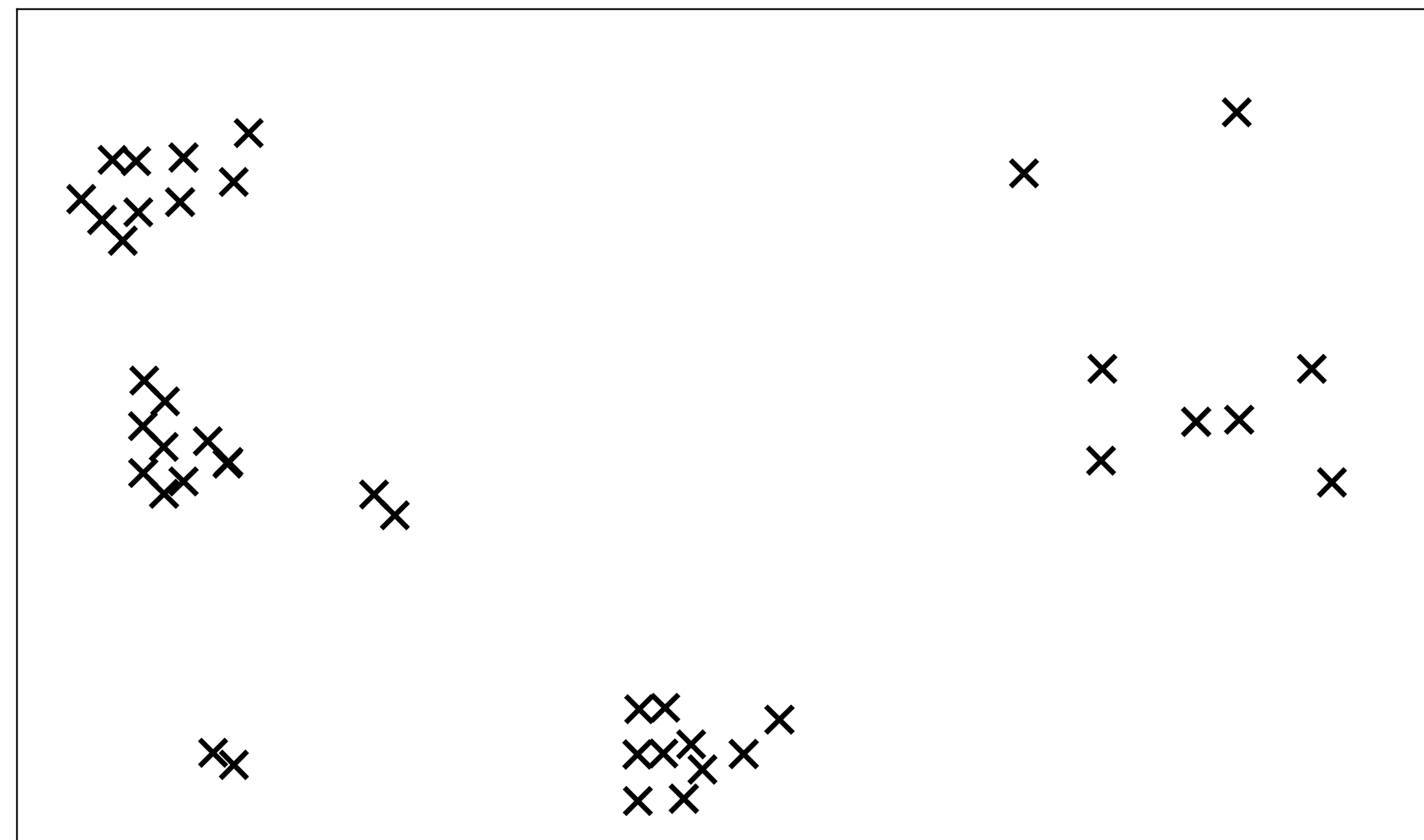
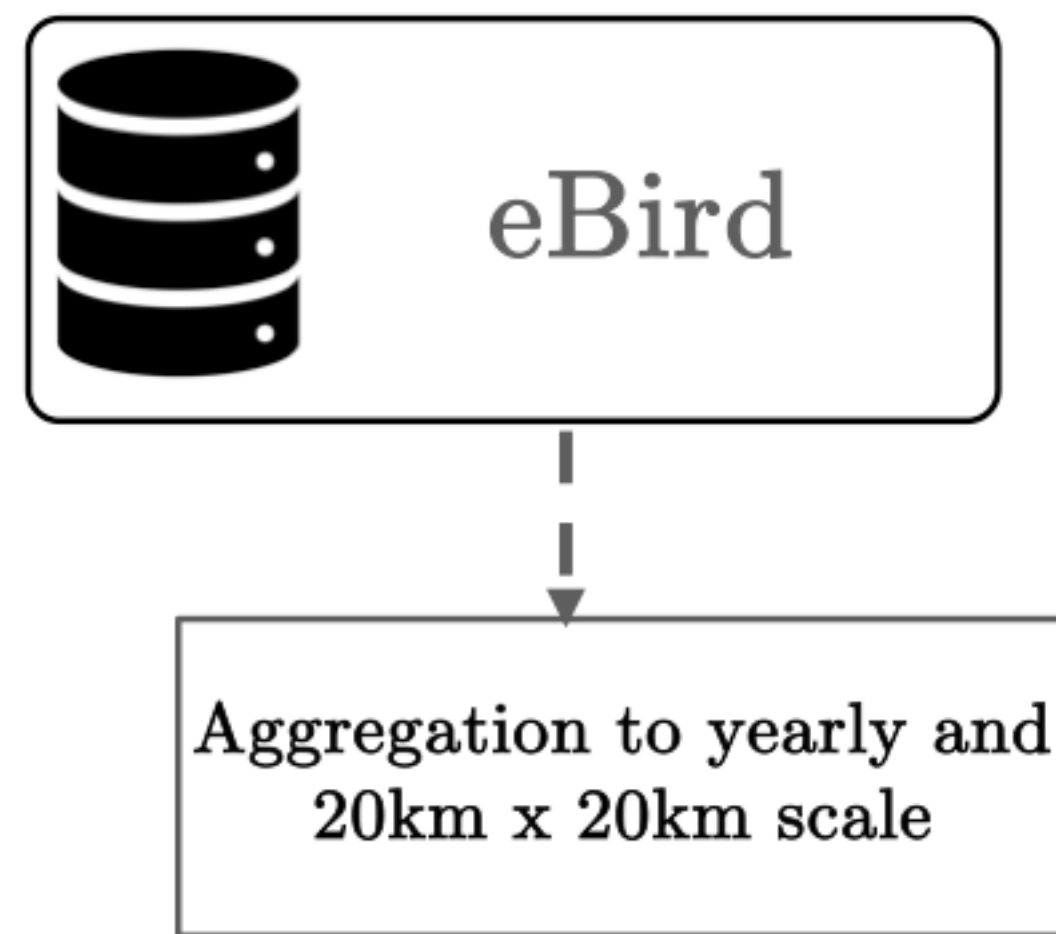
How does eBird work?

5. He finishes his checklist and submits his data to eBird
6. eBird internally verifies it, and it goes into their database



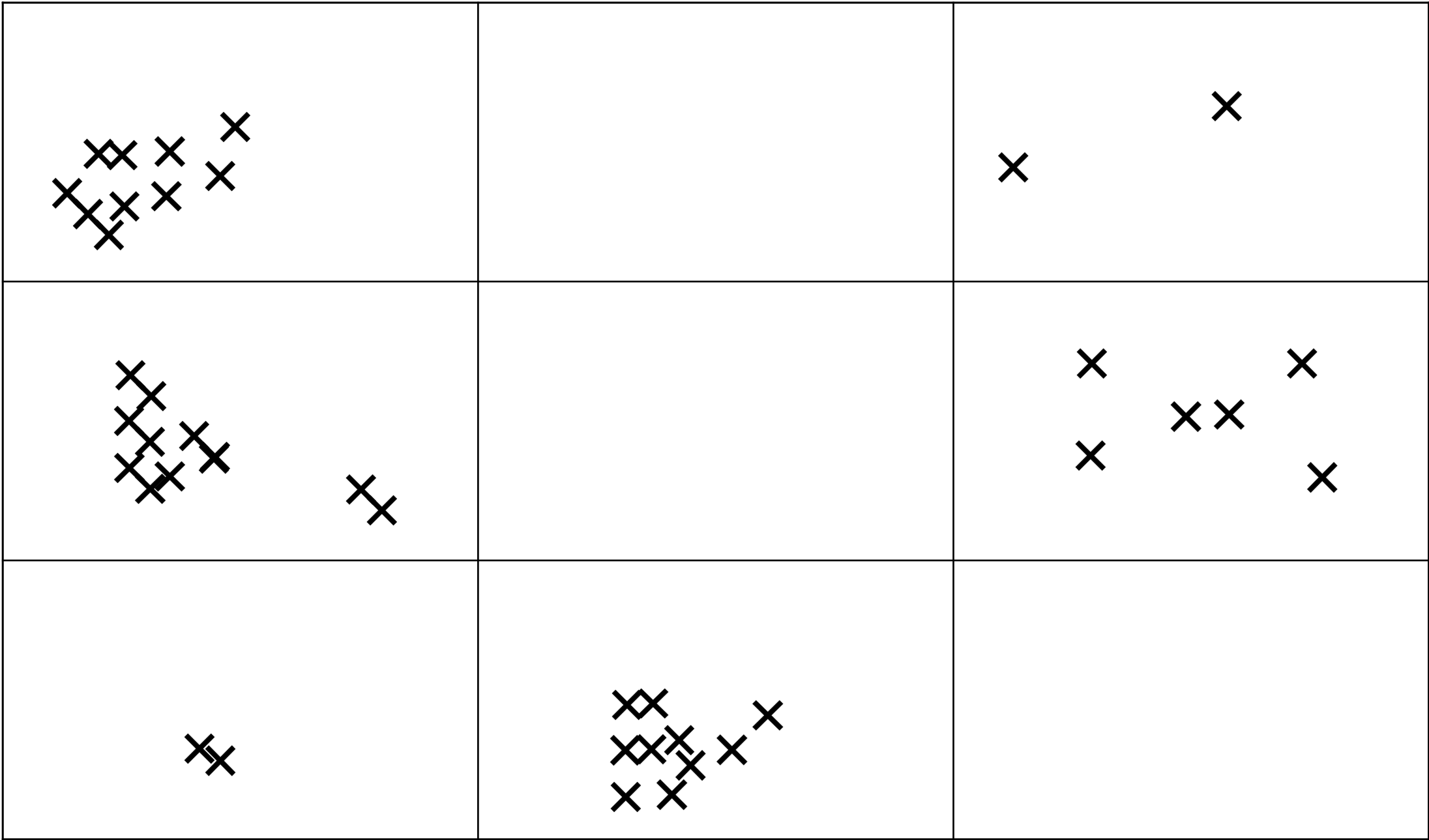
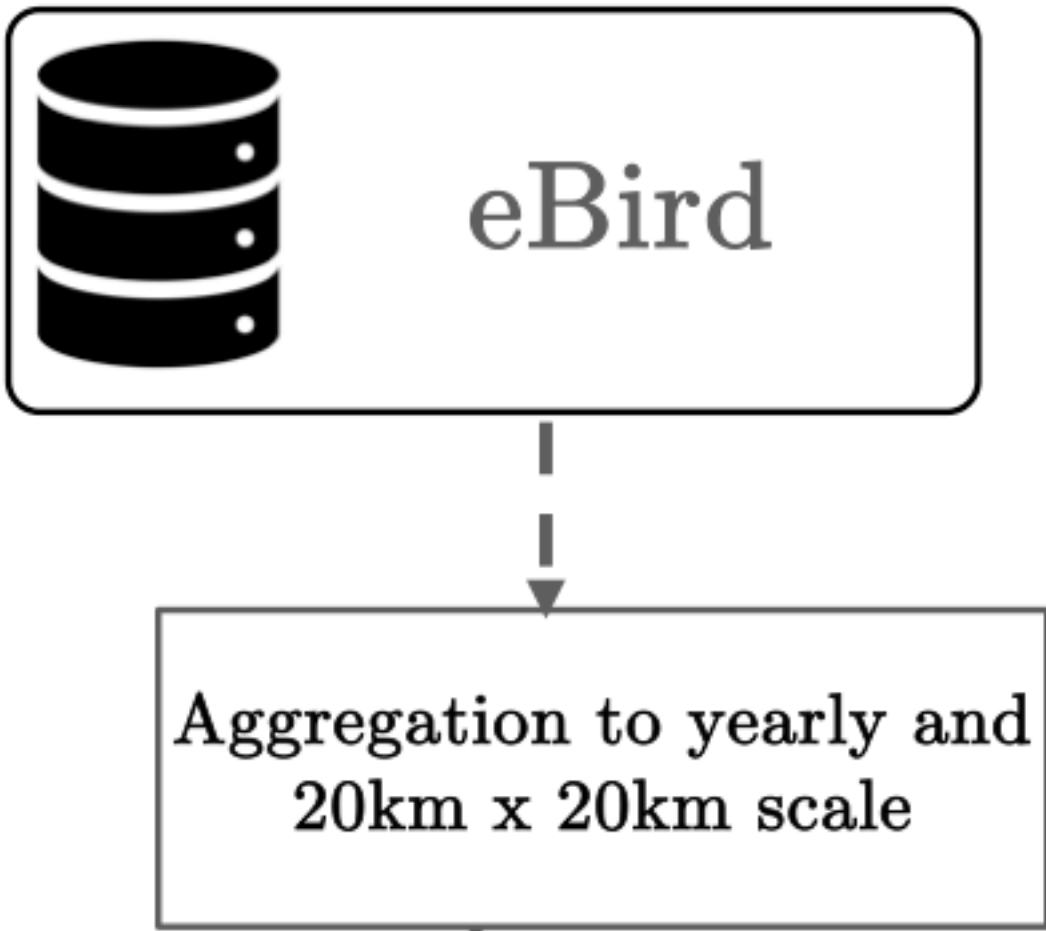
eBird data processing

Year 2020

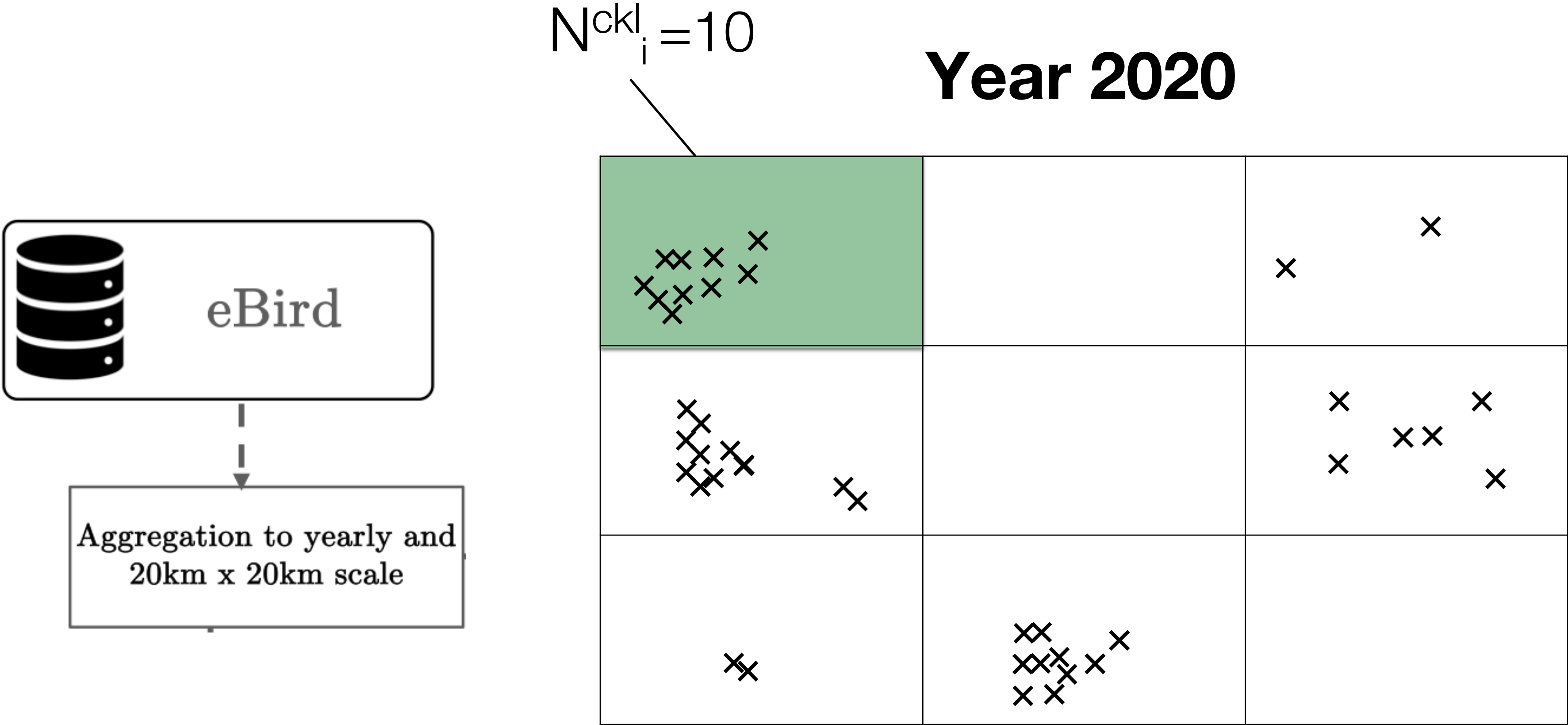


eBird data processing

Year 2020



eBird data processing



eBird data processing



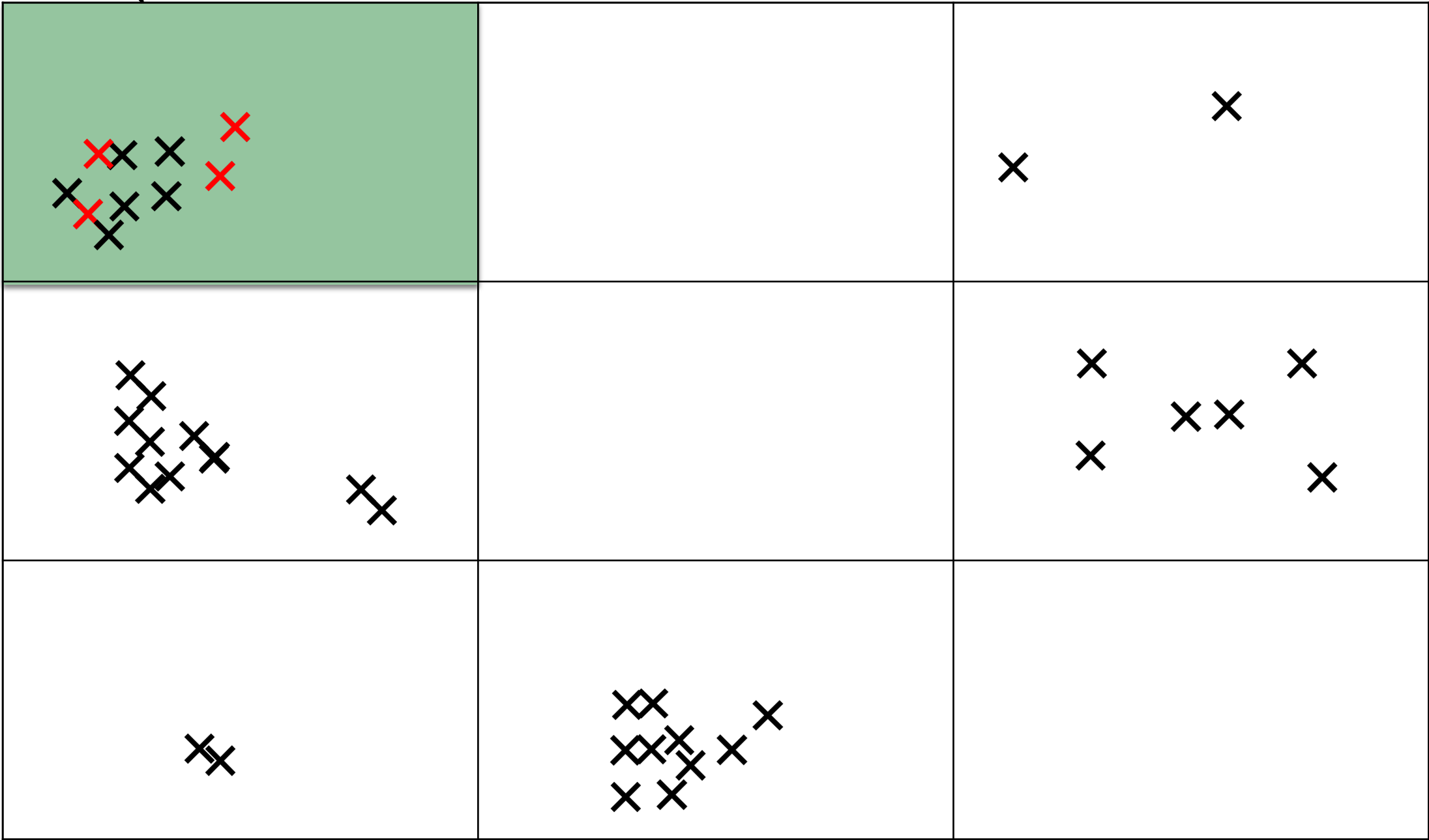
Purple Martin



Aggregation to yearly and
20km x 20km scale

$N^{spc}_i = 4$

Year 2020



eBird data processing

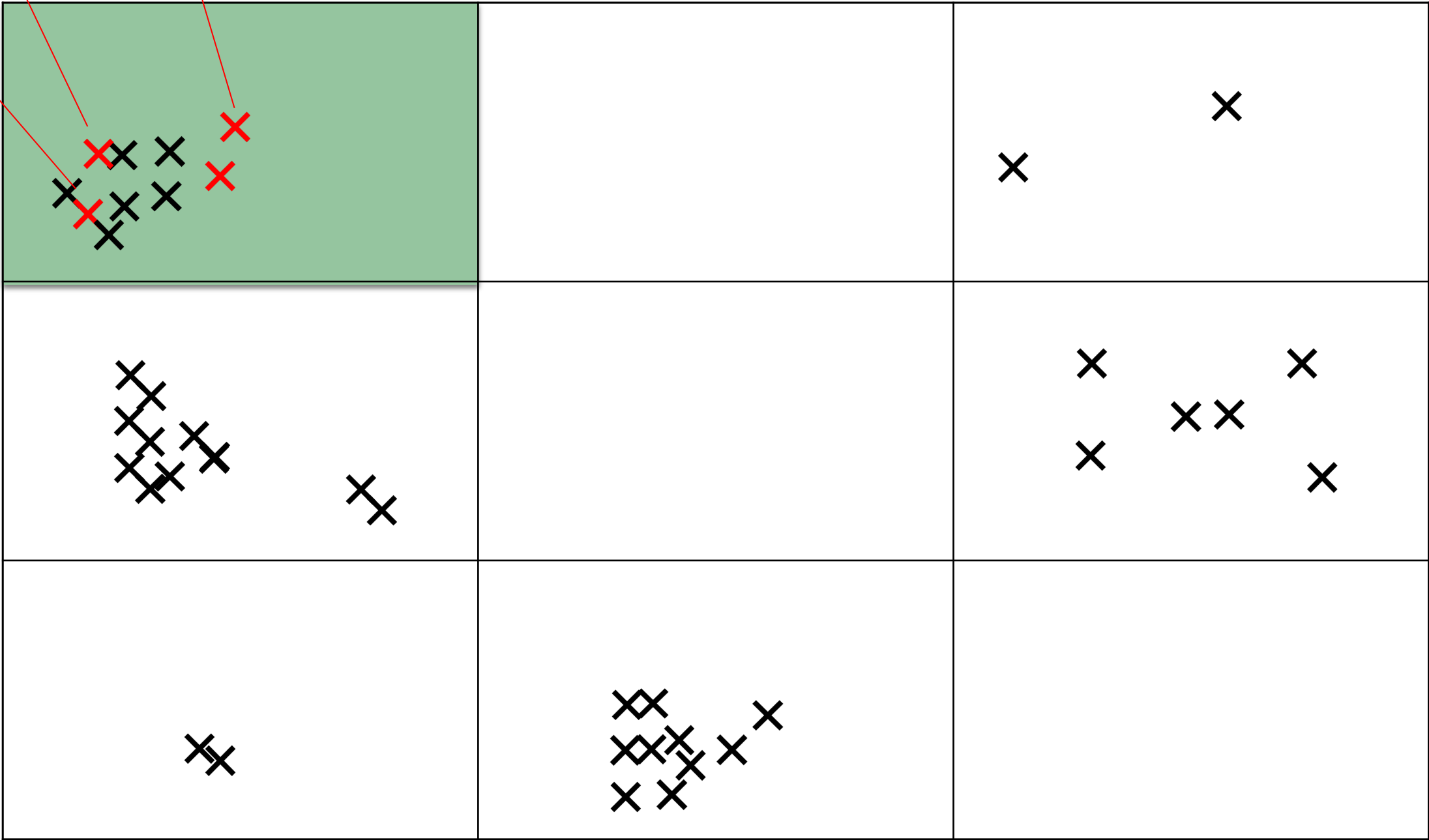


Purple Martin



Aggregation to yearly and 20km x 20km scale

$Y_{i,1}=28/03/2020$ $Y_{i,2}=02/04/2020$ $Y_{i,3}=05/04/2020$ **Year 2020**



eBird data processing



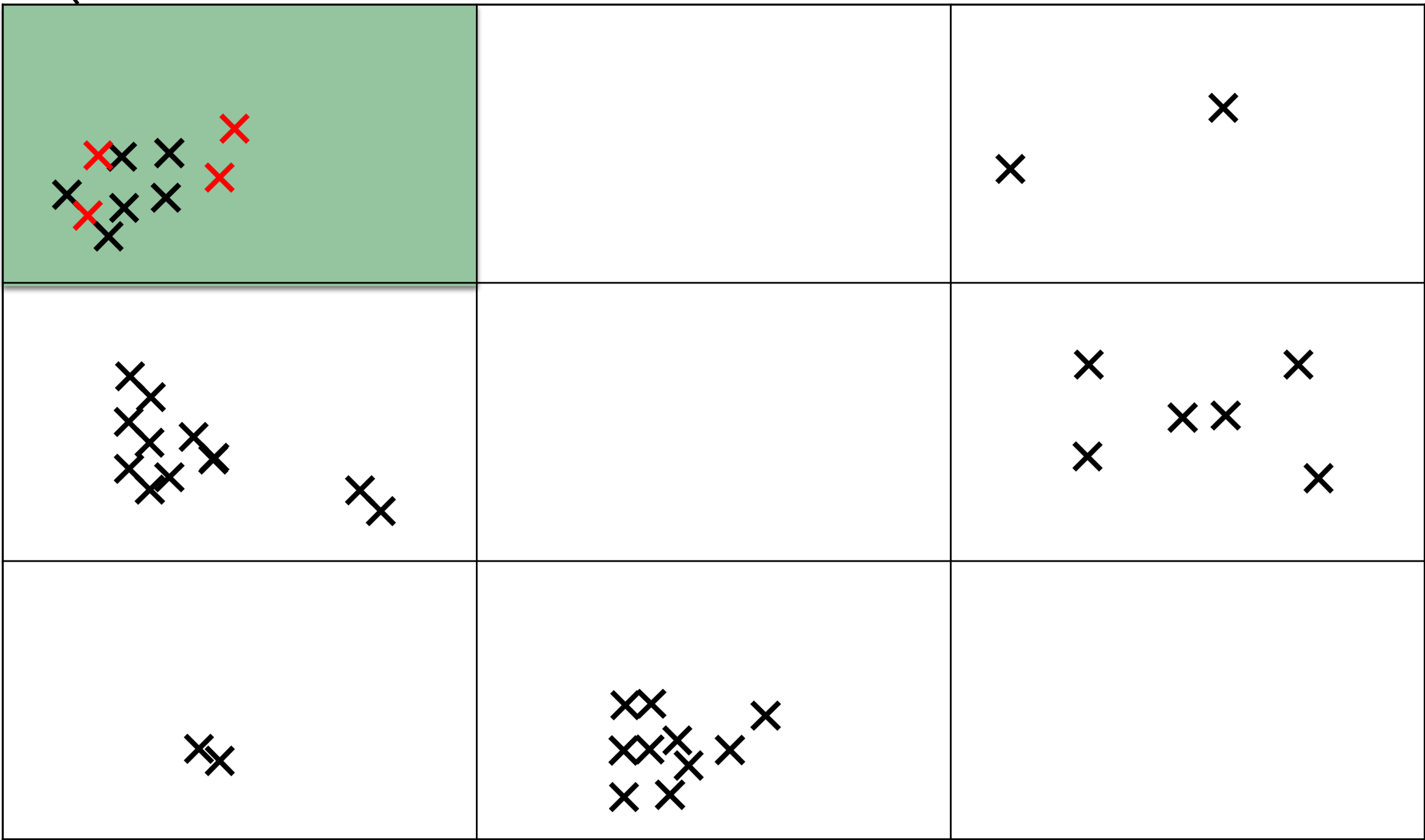
Purple Martin

Z'_i = First arrival date: 28/03/20

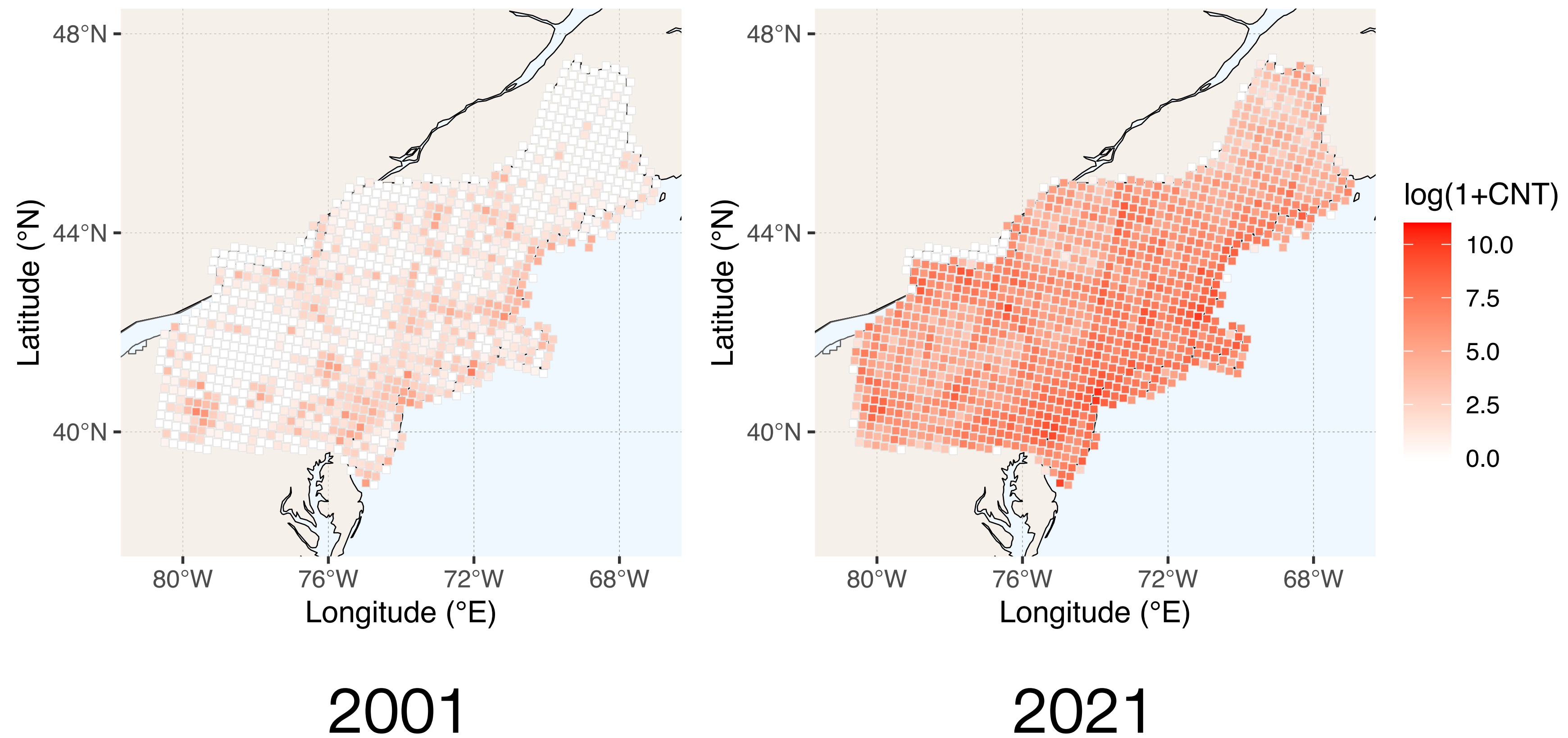
Year 2020



Aggregation to yearly and 20km x 20km scale



Strong temporal trends in reported occurrences

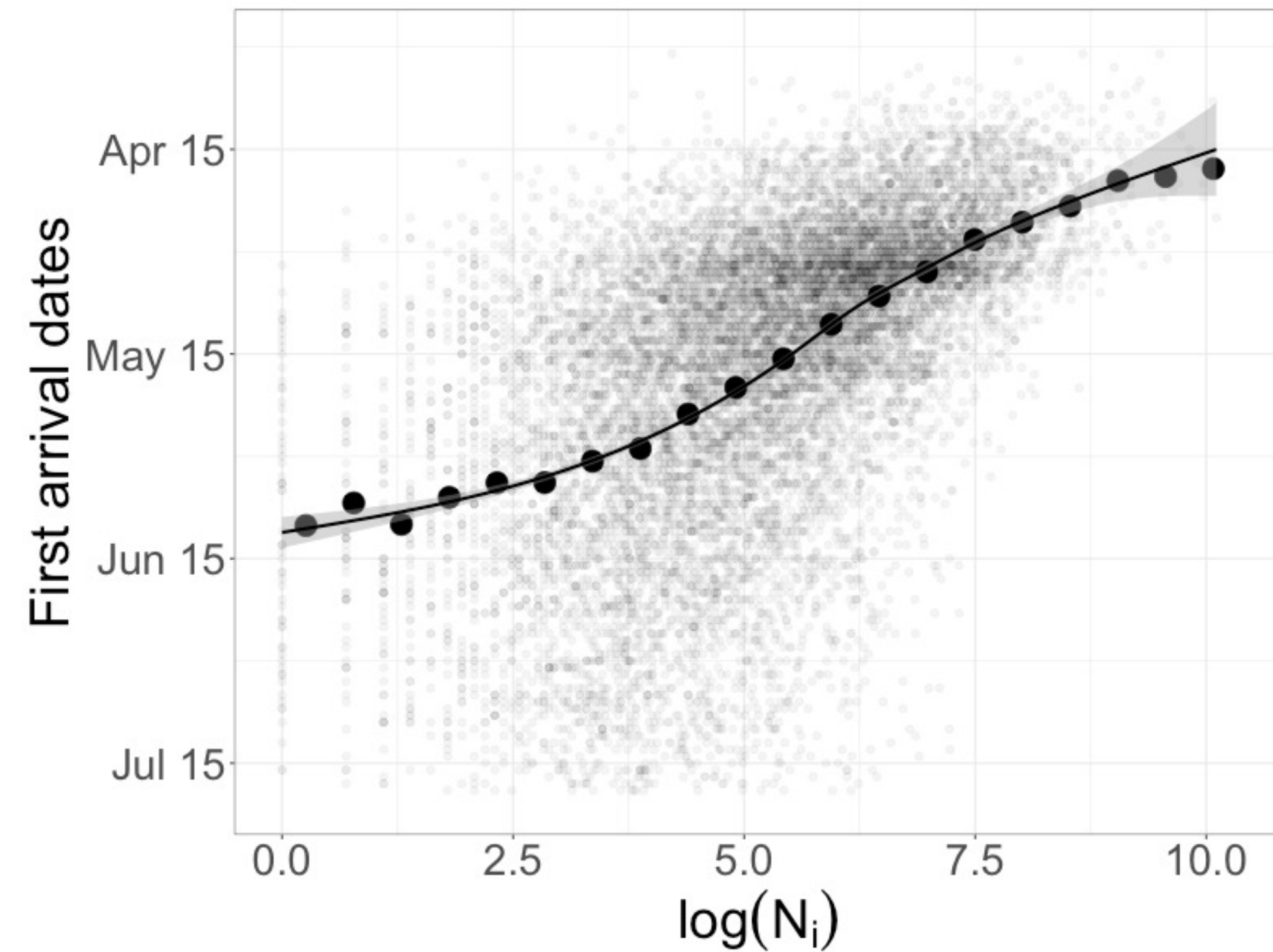


First arrival dates vs. checklist counts

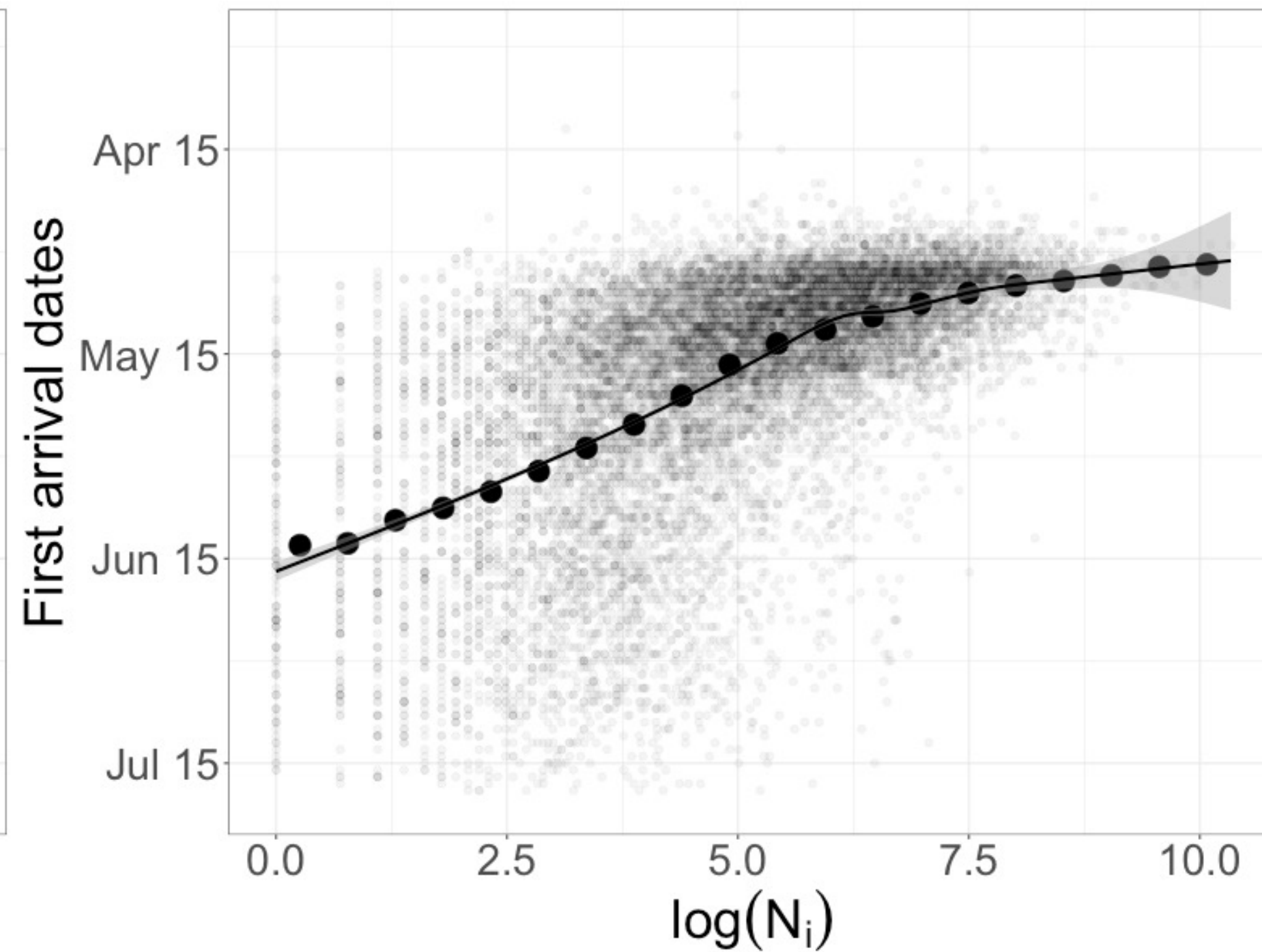
Citizen
Science

Bayesian
Hierarchical
Models

Chimney-Swift



**Chestnut-sided
Warbler**



“Observational effort”

Chimney-Swift

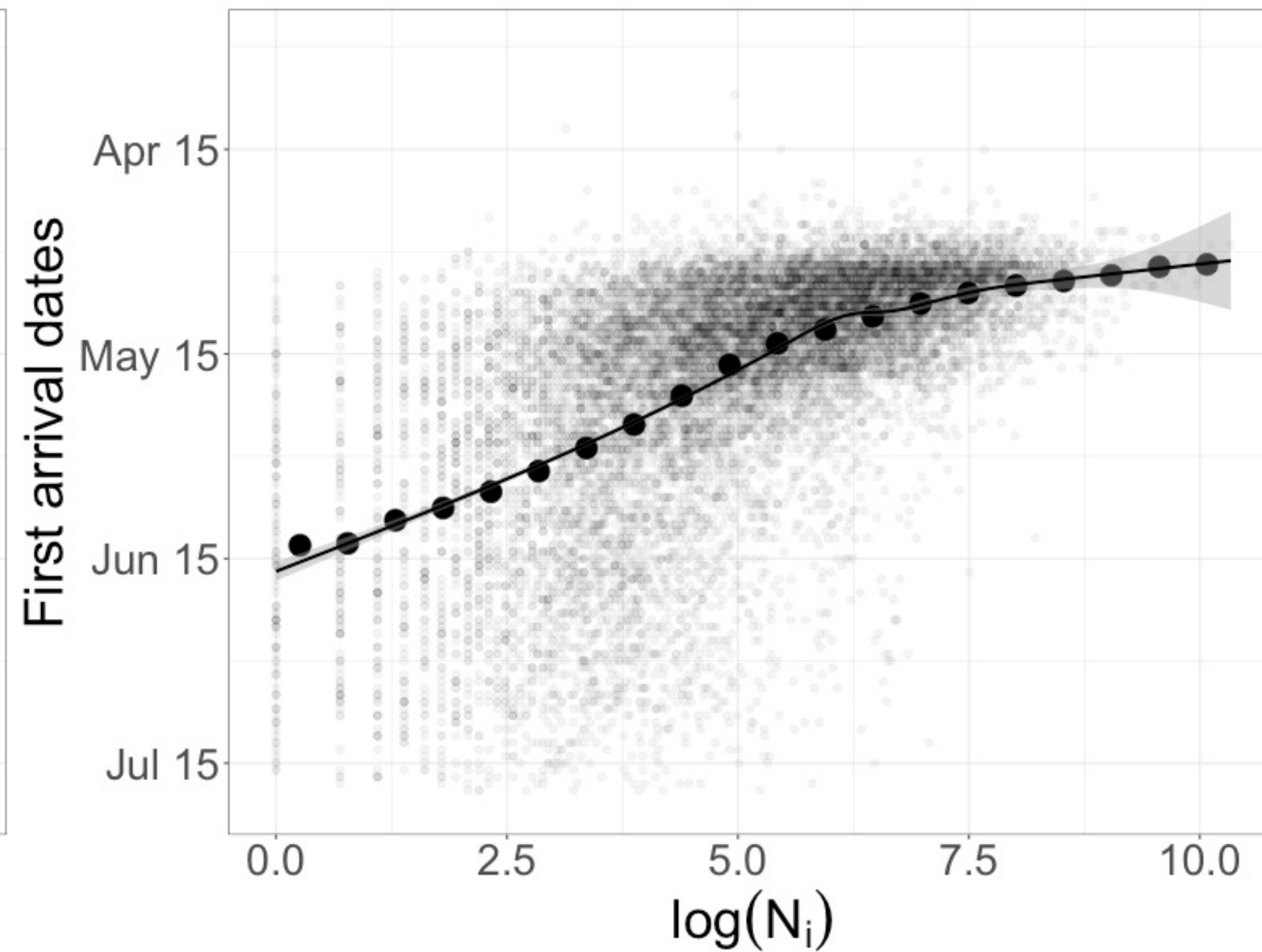
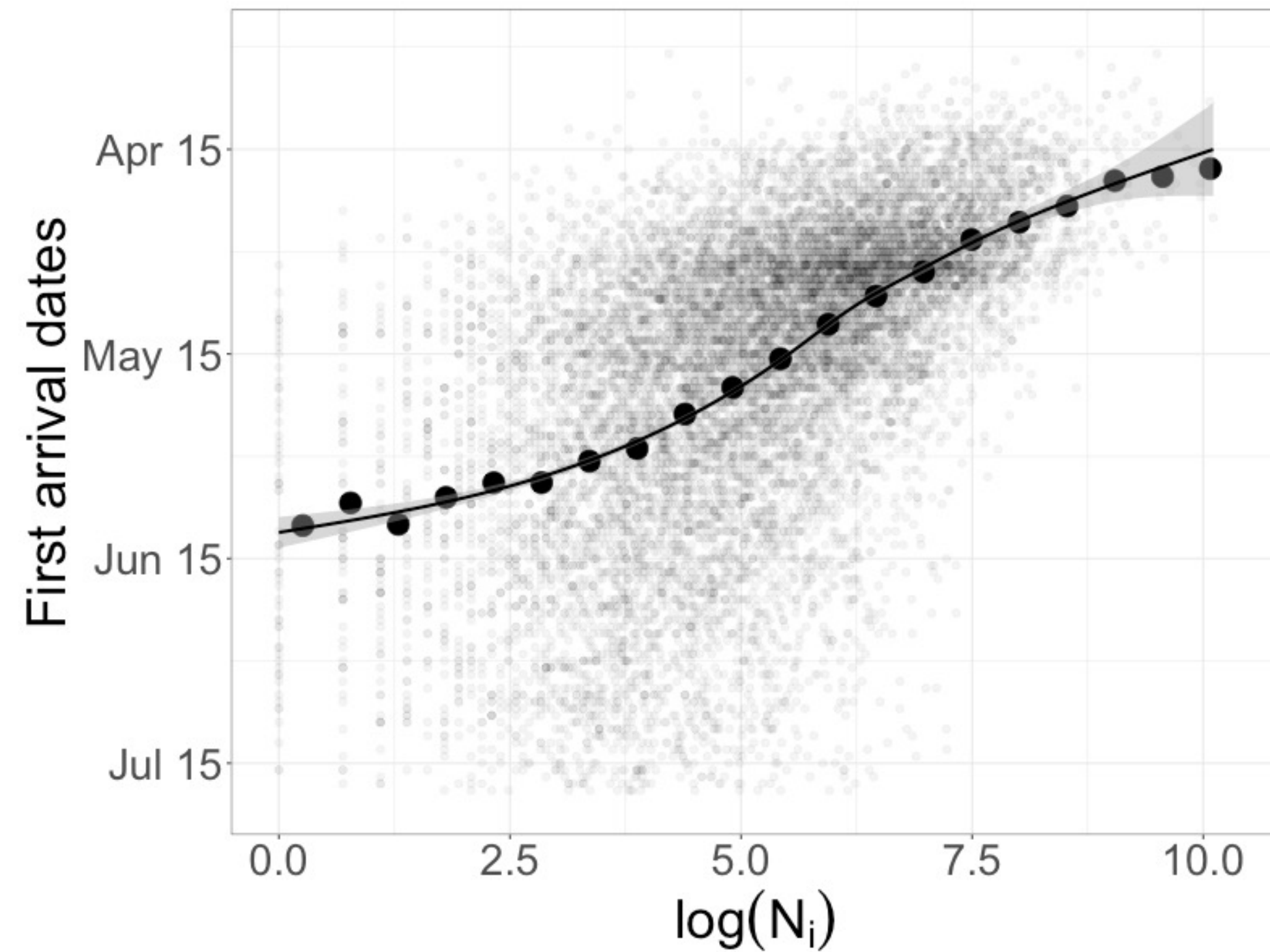


Chestnut-sided Warbler

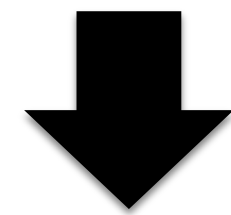


Citizen
Science

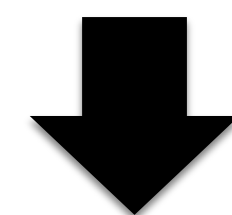
Bayesian
Hierarchical
Models



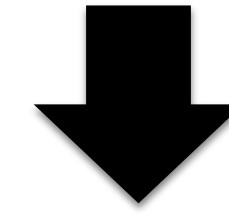
Observational effort = Preference + Activity



space-time
varying



captured by the
sampling intensity
for the checklists

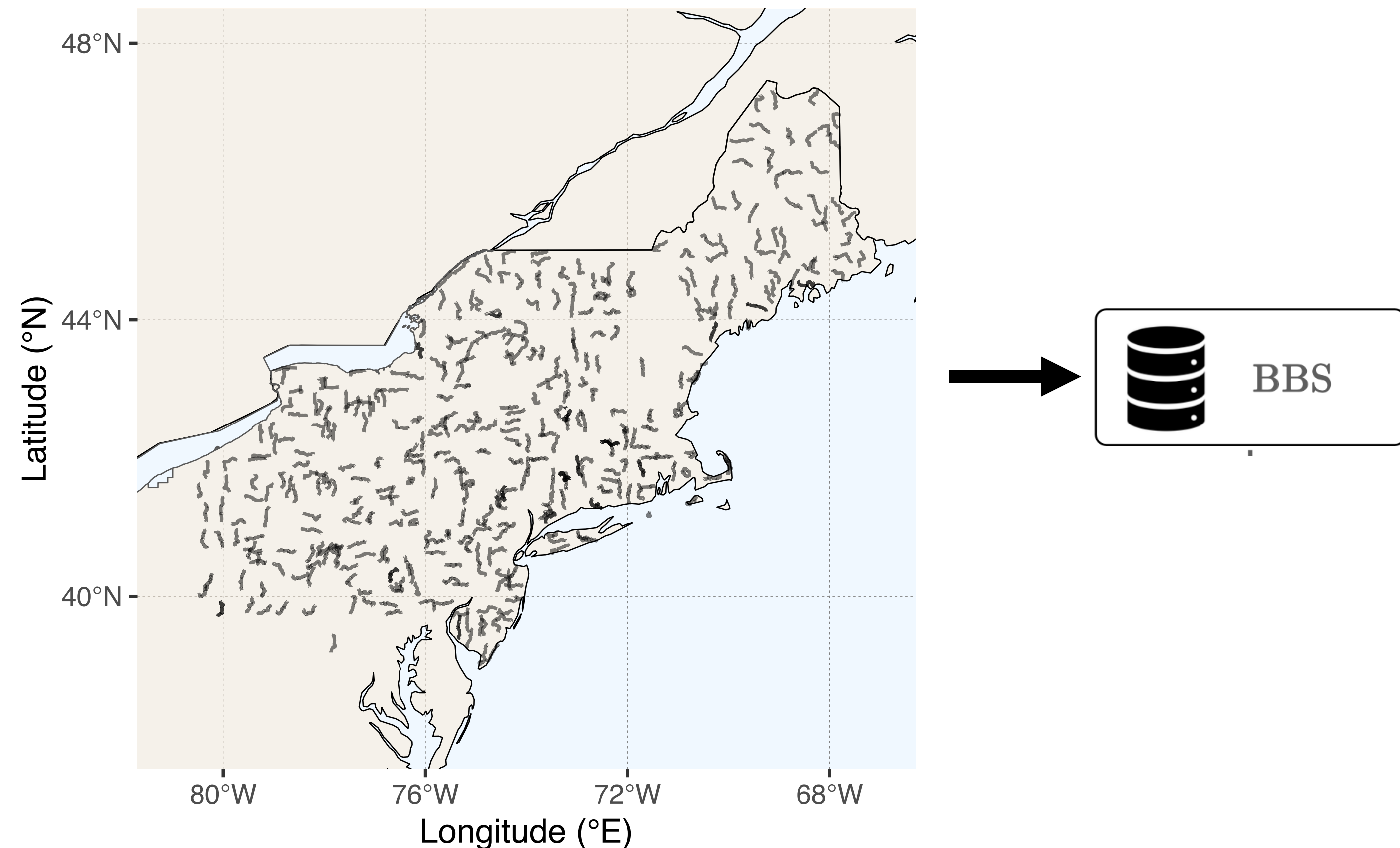


captured by the
(median) time spent
on the checklist

Breeding Bird Survey (BBS) sampling routes

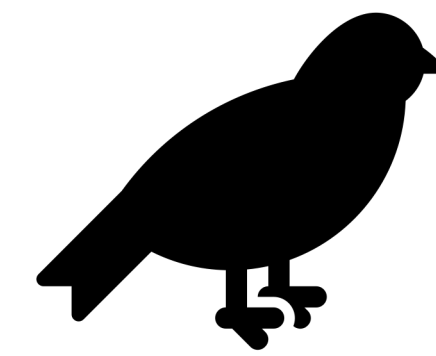
Citizen
Science

- For each route (~40km), bird occurrences are reported at 50 equidistant stops
- Complex data preprocessing (missing observations, missing stop coordinates, etc.)





MODEL



Modelling goals

- Fit a realistic model to first arrival data, conditional on covariates
- Correct for the observational bias from these datasets
- Use the model to make posterior predictions
- Interpolate spatially to locations not visited, in a reasonable way

A multi-response spatial regression system

Multi-response spatial regression

$$N_j^{\text{BBS}} \mid \lambda^{\text{BBS}}, \boldsymbol{\theta}_{\text{bbs}} \sim \text{Pois} \left\{ \sum_{k \in \text{route}_j} \omega_k \lambda^{\text{BBS}}(\mathbf{s}_k; \boldsymbol{\theta}_{\text{bbs}}) \right\},$$

$$N_i^{\text{ckl}} \mid \lambda^{\text{ckl}}, \boldsymbol{\theta}_{\text{ckl}} \sim \text{Pois} \{ \lambda^{\text{ckl}}(\mathbf{s}_i, t_i; \boldsymbol{\theta}_{\text{ckl}}) \},$$

$$N_i^{\text{spc}} \mid N_i^{\text{ckl}}, p^{\text{spc}}, \boldsymbol{\theta}_{\text{spc}} \sim \text{Bin} \{ N_i^{\text{ckl}}, p^{\text{spc}}(\mathbf{s}_i, t_i; \boldsymbol{\theta}_{\text{spc}}) \},$$

$$Z_i \mid \mu, \boldsymbol{\theta}_{\mu}, \sigma, \boldsymbol{\theta}_{\sigma} \sim \text{GEV} \{ \mu(\mathbf{s}_i, t_i; \boldsymbol{\theta}_{\mu}), \sigma(\mathbf{s}_i; \boldsymbol{\theta}_{\sigma}), \xi \},$$

where

$$\boldsymbol{\theta}_{\text{bbs}}, \boldsymbol{\theta}_{\text{ckl}}, \boldsymbol{\theta}_{\text{spc}}, \boldsymbol{\theta}_{\mu}, \boldsymbol{\theta}_{\sigma} \sim \text{Hyperpriors}$$

A multi-response spatial regression system

Multi-response spatial regression



$$\left\{ \begin{array}{l} N_j^{\text{BBS}} \mid \lambda^{\text{BBS}}, \theta_{\text{bbs}} \sim \text{Pois} \left\{ \sum_{k \in \text{route}_j} \omega_k \lambda^{\text{BBS}}(\mathbf{s}_k; \theta_{\text{bbs}}) \right\}, \end{array} \right.$$



$$\left\{ \begin{array}{l} N_i^{\text{ckl}} \mid \lambda^{\text{ckl}}, \theta_{\text{ckl}} \sim \text{Pois} \{ \lambda^{\text{ckl}}(\mathbf{s}_i, t_i; \theta_{\text{ckl}}) \}, \\ N_i^{\text{spc}} \mid N_i^{\text{ckl}}, p^{\text{spc}}, \theta_{\text{spc}} \sim \text{Bin} \{ N_i^{\text{ckl}}, p^{\text{spc}}(\mathbf{s}_i, t_i; \theta_{\text{spc}}) \}, \\ Z_i \mid \mu, \theta_{\mu}, \sigma, \theta_{\sigma} \sim \text{GEV} \{ \mu(\mathbf{s}_i, t_i; \theta_{\mu}), \sigma(\mathbf{s}_i; \theta_{\sigma}), \xi \}, \end{array} \right.$$

where

$$\theta_{\text{bbs}}, \theta_{\text{ckl}}, \theta_{\text{spc}}, \theta_{\mu}, \theta_{\sigma} \sim \text{Hyperpriors}$$

Sharing random effects

Multi-response spatial regression



$$\begin{aligned}
 & \left\{ \begin{aligned} & N_j^{\text{BBS}} \mid \lambda^{\text{BBS}}, \theta_{\text{bbs}} \sim \text{Pois} \left\{ \sum_{k \in \text{route}_j} \omega_k \lambda^{\text{BBS}}(\mathbf{s}_k; \theta_{\text{bbs}}) \right\}, \\ & N_i^{\text{ckl}} \mid \lambda^{\text{ckl}}, \theta_{\text{ckl}} \sim \text{Pois} \{ \lambda^{\text{ckl}}(\mathbf{s}_i, t_i; \theta_{\text{ckl}}) \}, \\ & N_i^{\text{spc}} \mid N_i^{\text{ckl}}, p^{\text{spc}}, \theta_{\text{spc}} \sim \text{Bin} \{ N_i^{\text{ckl}}, p^{\text{spc}}(\mathbf{s}_i, t_i; \theta_{\text{spc}}) \}, \\ & Z_i \mid \mu, \theta_{\mu}, \sigma, \theta_{\sigma} \sim \text{GEV} \{ \mu(\mathbf{s}_i, t_i; \theta_{\mu}), \sigma(\mathbf{s}_i; \theta_{\sigma}), \xi \}, \end{aligned} \right.
 \end{aligned}$$

Diagram illustrating the sharing of random effects between BBS and eBird data. The BBS data is modeled by $N_j^{\text{BBS}} \mid \lambda^{\text{BBS}}, \theta_{\text{bbs}} \sim \text{Pois} \left\{ \sum_{k \in \text{route}_j} \omega_k \lambda^{\text{BBS}}(\mathbf{s}_k; \theta_{\text{bbs}}) \right\}$. The eBird data is modeled by $N_i^{\text{ckl}} \mid \lambda^{\text{ckl}}, \theta_{\text{ckl}} \sim \text{Pois} \{ \lambda^{\text{ckl}}(\mathbf{s}_i, t_i; \theta_{\text{ckl}}) \}$, $N_i^{\text{spc}} \mid N_i^{\text{ckl}}, p^{\text{spc}}, \theta_{\text{spc}} \sim \text{Bin} \{ N_i^{\text{ckl}}, p^{\text{spc}}(\mathbf{s}_i, t_i; \theta_{\text{spc}}) \}$, and $Z_i \mid \mu, \theta_{\mu}, \sigma, \theta_{\sigma} \sim \text{GEV} \{ \mu(\mathbf{s}_i, t_i; \theta_{\mu}), \sigma(\mathbf{s}_i; \theta_{\sigma}), \xi \}$. The spatial process $X^{\text{niche}}(\cdot) \sim \mathcal{GP}(\omega_2)$ is shared between the two datasets, influencing the BBS data through ω_k and the eBird data through $\mu(\mathbf{s}_i, t_i; \theta_{\mu})$ and $\sigma(\mathbf{s}_i; \theta_{\sigma})$.

where

$$\theta_{\text{bbs}}, \theta_{\text{ckl}}, \theta_{\text{spc}}, \theta_{\mu}, \theta_{\sigma} \sim \text{Hyperpriors}$$

Sharing random effects

Multi-response spatial regression



$$\begin{aligned}
 & \left\{ \begin{aligned} N_j^{\text{BBS}} \mid \lambda^{\text{BBS}}, \theta_{\text{bbs}} &\sim \text{Pois} \left\{ \sum_{k \in \text{route}_j} \omega_k \lambda^{\text{BBS}}(\mathbf{s}_k; \theta_{\text{bbs}}) \right\}, \\ N_i^{\text{ckl}} \mid \lambda^{\text{ckl}}, \theta_{\text{ckl}} &\sim \text{Pois} \{ \lambda^{\text{ckl}}(\mathbf{s}_i, t_i; \theta_{\text{ckl}}) \}, \\ N_i^{\text{spc}} \mid N_i^{\text{ckl}}, p^{\text{spc}}, \theta_{\text{spc}} &\sim \text{Bin} \{ N_i^{\text{ckl}}, p^{\text{spc}}(\mathbf{s}_i, t_i; \theta_{\text{spc}}) \}, \\ Z_i \mid \mu, \theta_{\mu}, \sigma, \theta_{\sigma} &\sim \text{GEV} \{ \mu(\mathbf{s}_i, t_i; \theta_{\mu}), \sigma(\mathbf{s}_i; \theta_{\sigma}), \xi \}, \end{aligned} \right.
 \end{aligned}$$

$X^{\text{pref}}(\cdot) \sim \mathcal{GP}(\omega_1)$

where

$$\theta_{\text{bbs}}, \theta_{\text{ckl}}, \theta_{\text{spc}}, \theta_{\mu}, \theta_{\sigma} \sim \text{Hyperpriors}$$

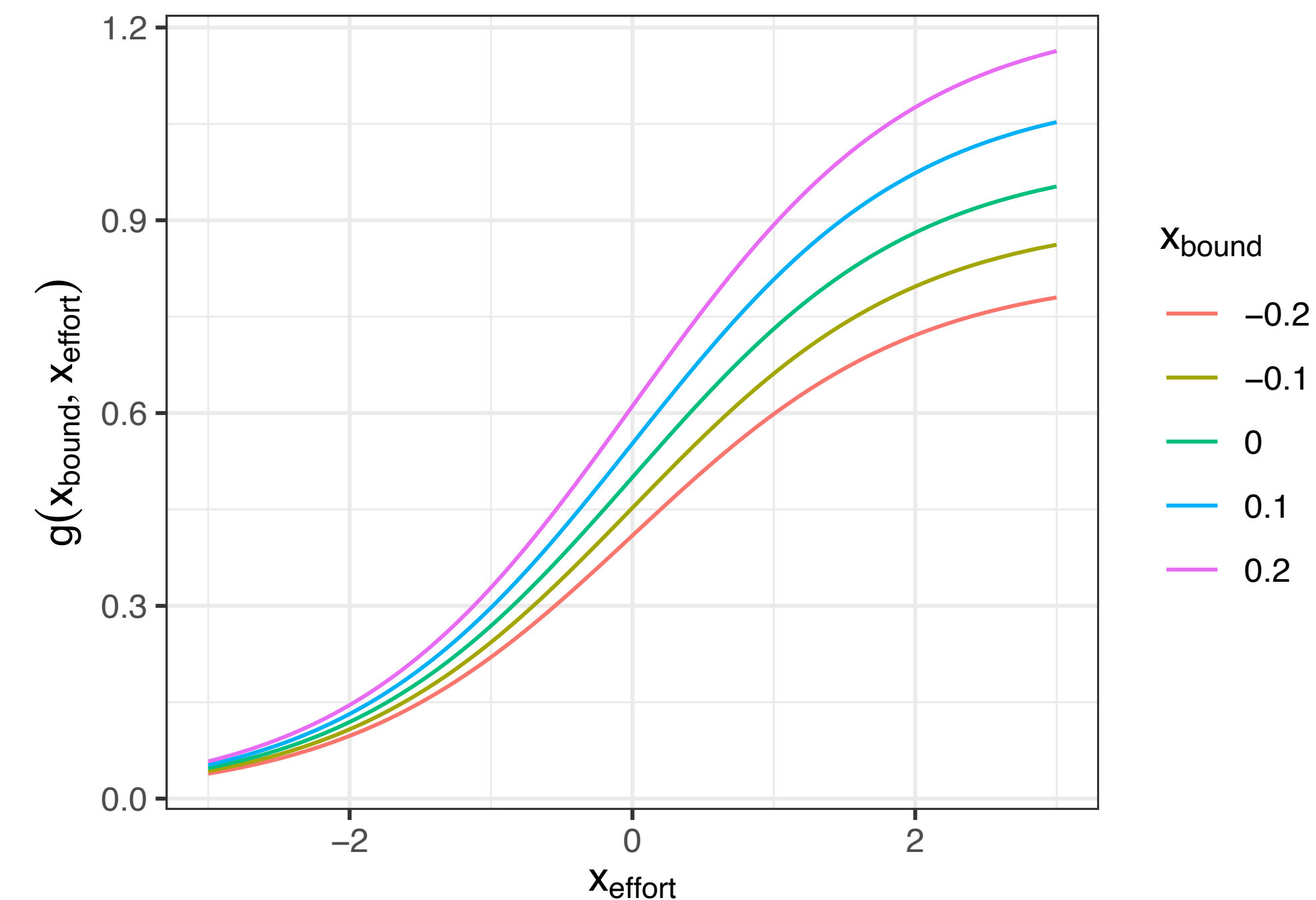
Saturating effect of observational effort

Bayesian
Hierarchical
Models

Extreme-
Value
Theory

- Observed first arrival is biased towards later dates for low effort but is the true one for very high effort
- Implementation: $Z_i \sim \text{GEV}(\mu_i, \sigma_i)$ with $\mu_i = g(\text{Predictors}_i, \text{Effort}_i)$
 - Nonlinear function g reaches (unknown) finite upper bound for very high effort
 - Infer g from data
 - Set very high effort for bias-corrected predictions

⚠ Source of high computational complexity



Goodness-of-fit of estimated models

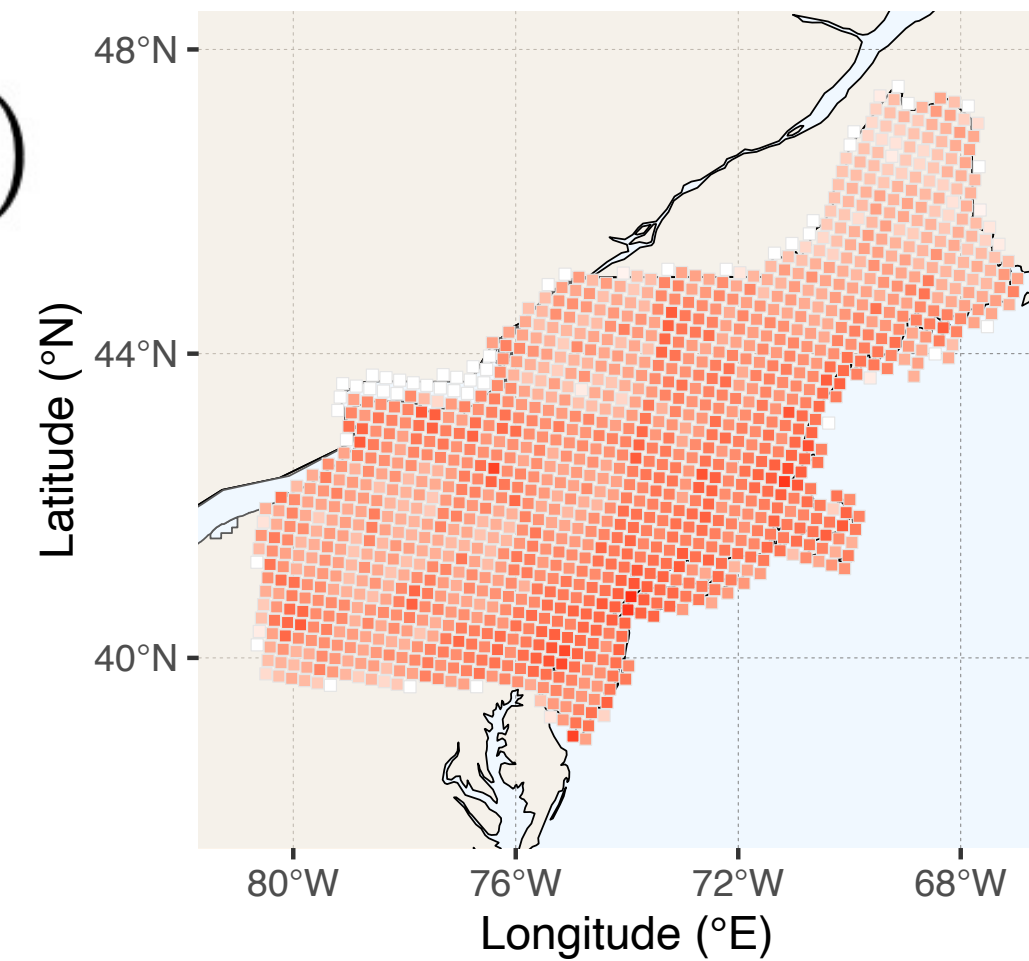
- Generally good match of eBird observations (left maps) with posterior means (right maps)
- Slight differences due to information shared from BBS

Example species:

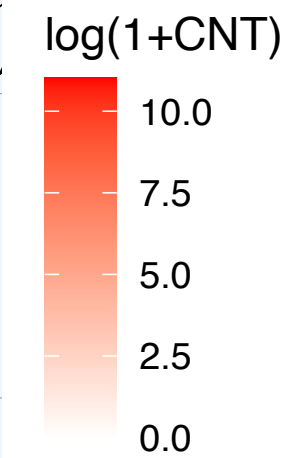
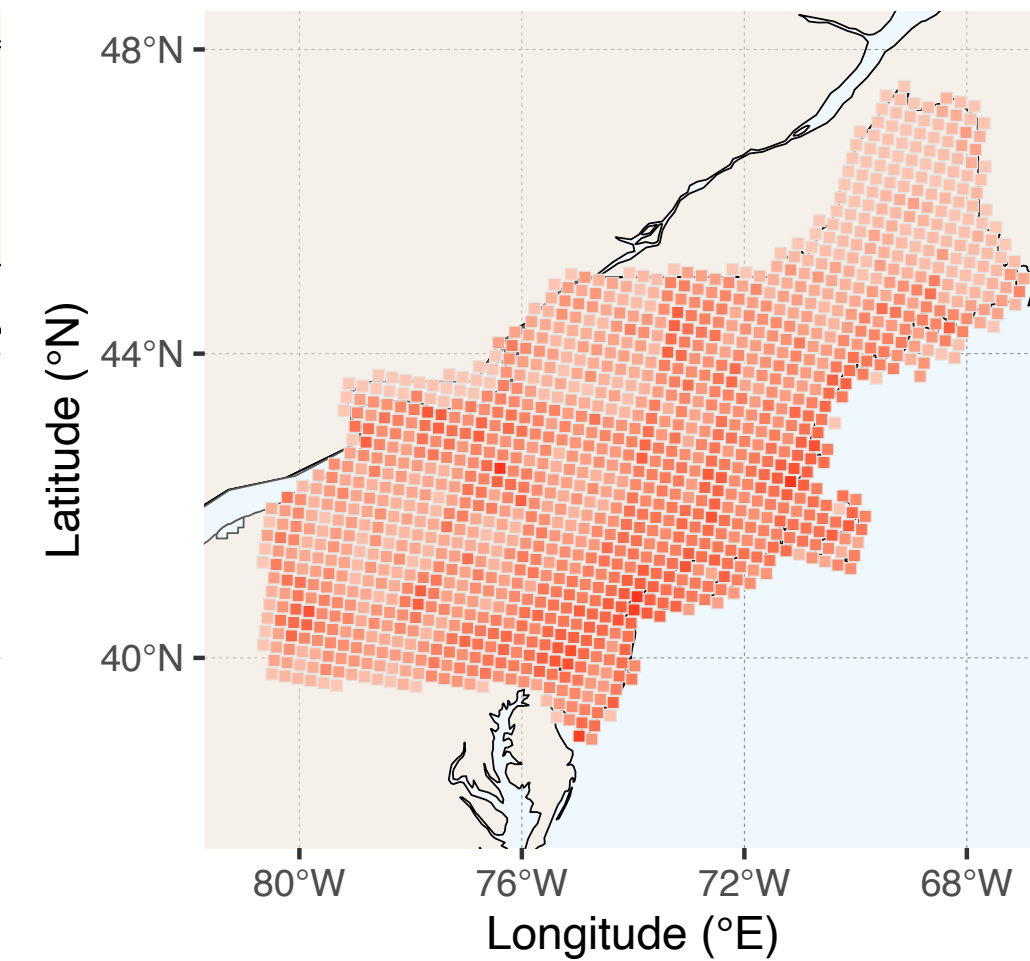


Great Crested Flycatcher

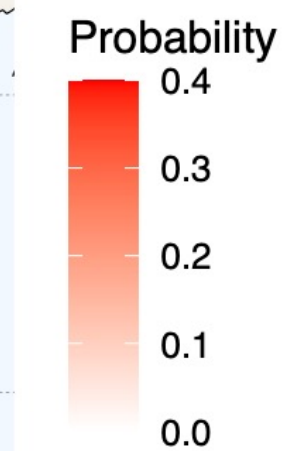
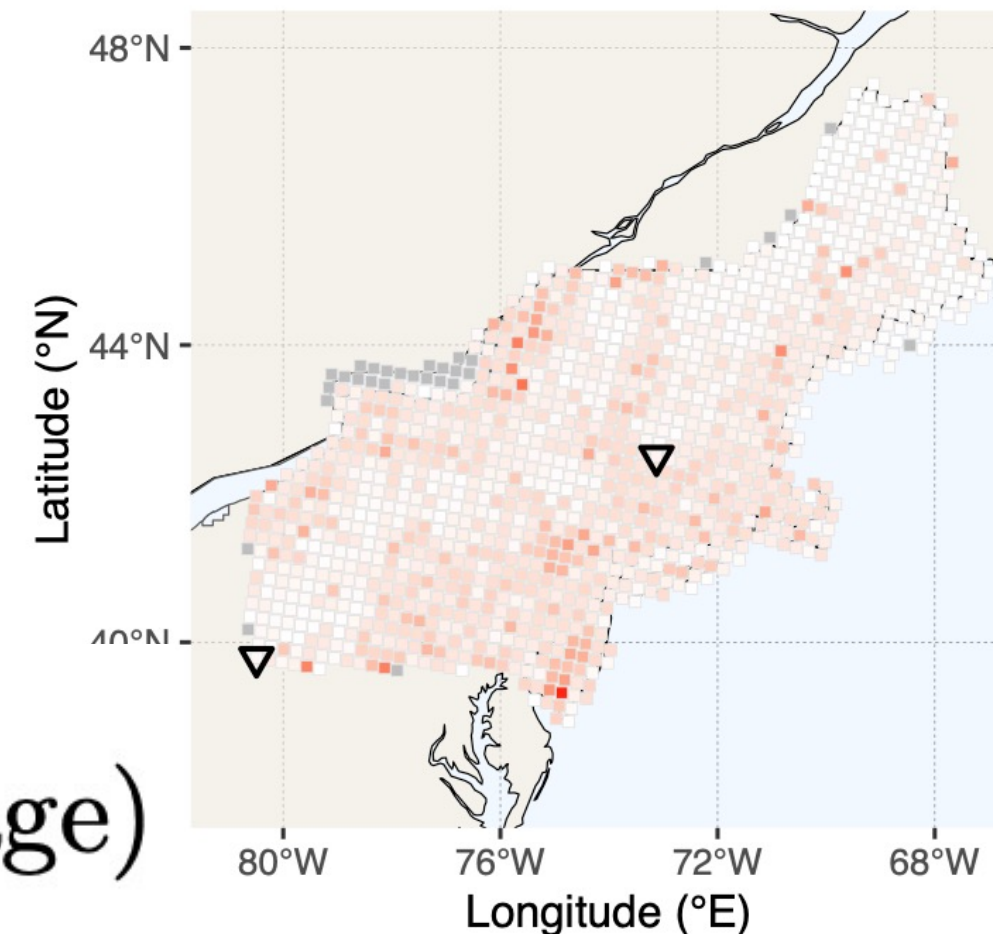
$$N_i^{\text{ckl}}(2021)$$



$$\hat{\lambda}_i^{\text{ckl}}(2021)$$



$$\frac{N_i^{\text{spc}}}{N_i^{\text{ckl}}}(\text{Time-average})$$



$$\hat{p}_i^{\text{spc}}(\text{Time-average})$$

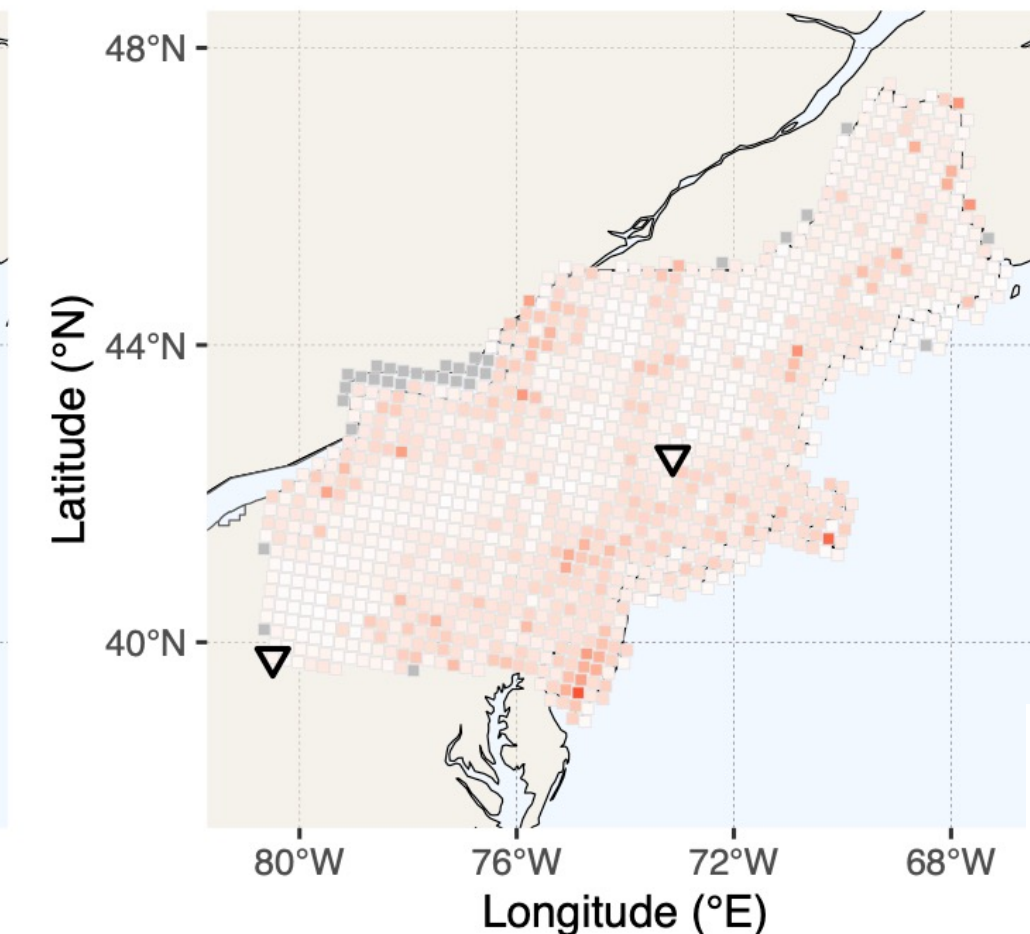


Illustration of bias-corrected prediction of first arrivals (2022)

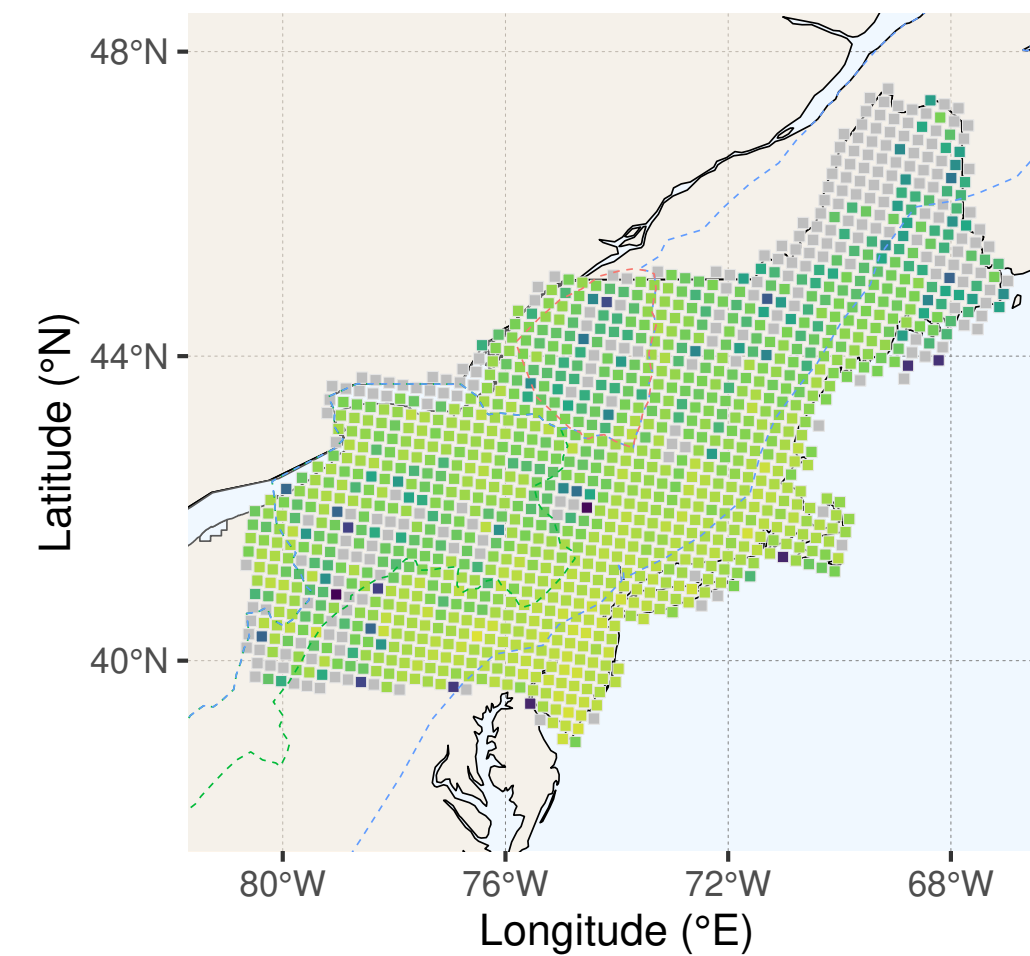
- Based on Generalised Extreme-Value response
- Bias-corrected prediction by fixing saturated observational effort

$$Z_i \mid \mu, \theta^\mu, \sigma, \theta_\sigma \sim \text{GEV}\{\mu(\mathbf{s}_i, t_i; \theta_\mu), \sigma(\mathbf{s}_i; \theta_\sigma), \xi\}$$

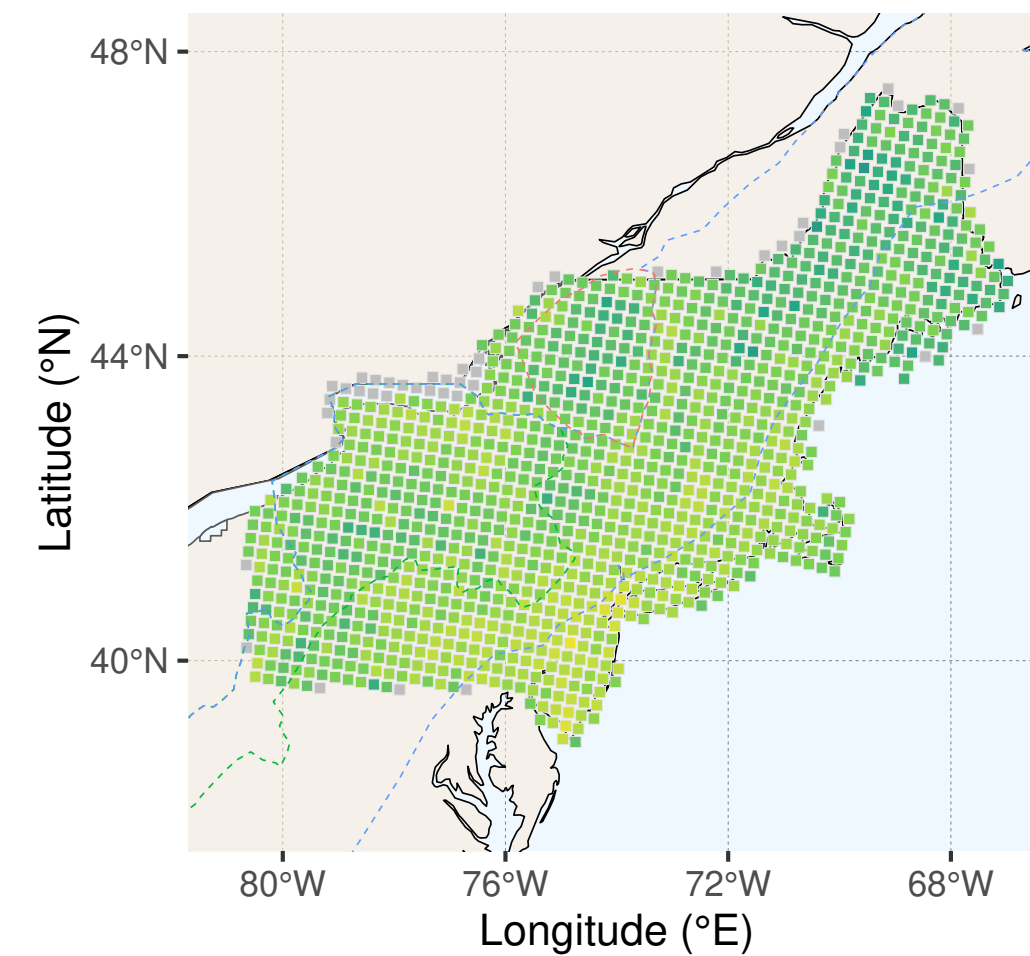


Great Crested
Flycatcher

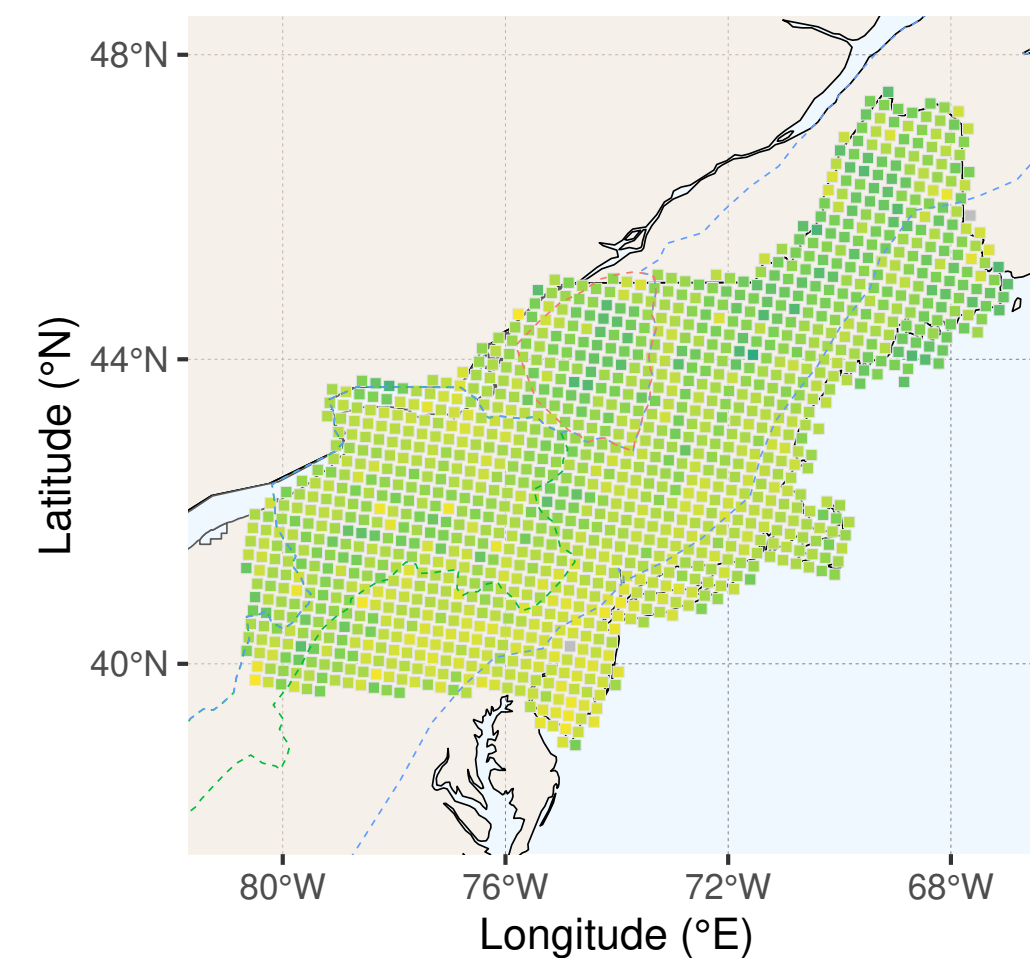
Observed



Posterior



Posterior predictive
→ Saturated effort



Posterior predictive
→ Saturated effort
→ Species niche

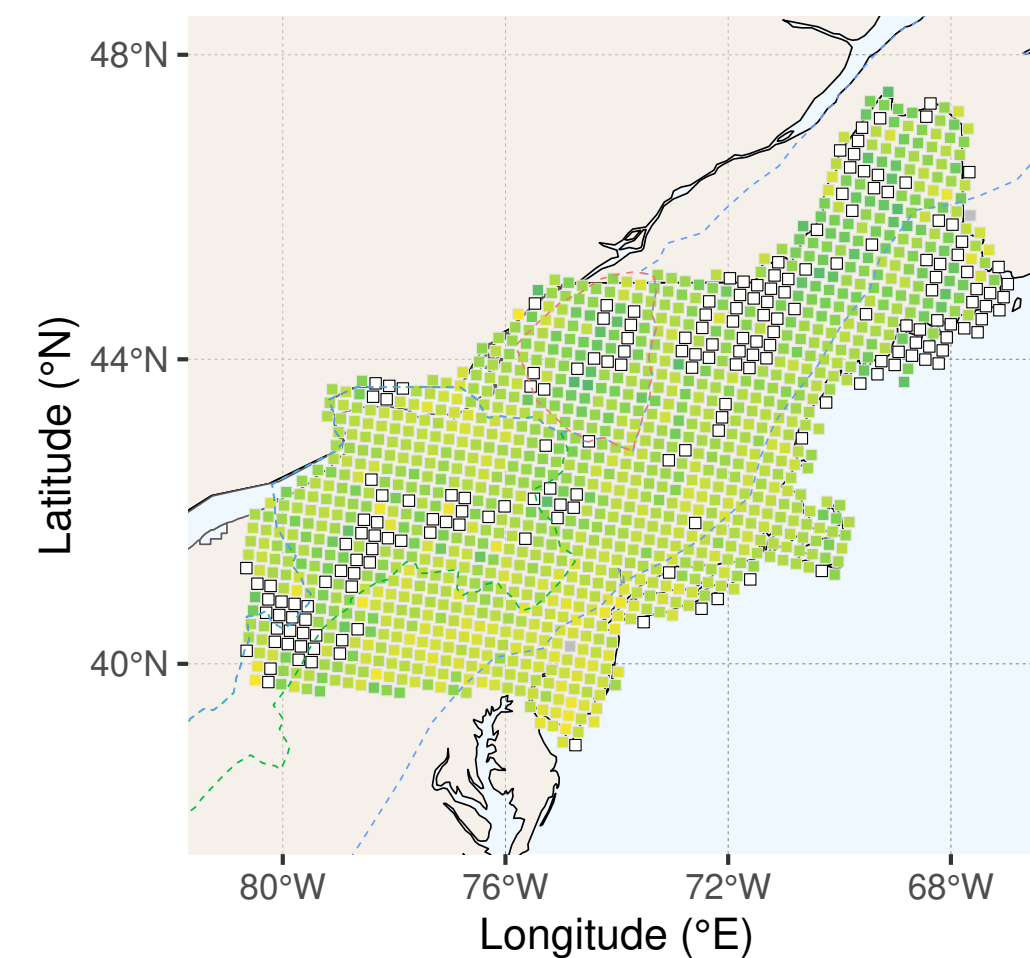


Illustration of bias-corrected prediction of first arrivals (2022), cont'd

- Table of estimated key parameters and first arrival dates for two pixels
- Estimated (not bias-corrected) first arrivals tend to occur relatively earlier for
 - higher Preference,
 - higher Activity and
 - in the core area of the niche



Species	Chimney Swift	Great Crested Flycatcher	Chestnut-sided Warbler	Purple Martin
$\hat{\theta}^{\text{pref}}$	0.191 (0.184,0.202)	0.204 (0.199,0.21)	0.187 (0.183,0.191)	0.2 (0.178,0.217)
$\hat{\theta}^{\text{act}}$	-0.15 (-0.217,-0.061)	-0.818 (-0.911,-0.696)	-0.548 (-0.619,-0.454)	-0.03 (-0.269,0.236)
$\hat{\theta}^{\text{niche-GEV}} (\times 10^{-2})$	4.9 (4.664,5.134)	4 (3.894,4.133)	0.2 (0.17,0.278)	6 (5.541,6.443)
Observed	NA	NA	NA	NA
Predicted	09/05	03/05	21/05	07/06
Debiased	03/04	13/04	03/05	28/03
Observed	01/05	04/05	04/05	29/06
Predicted	09/05	15/05	12/05	12/05
Debiased	22/04	05/05	03/05	07/04

Discussion: Ecological data fusion using latent processes

Bayesian Hierarchical Models

- Incomplete and biased observation of true processes
- Interpretable latent processes for effort and relevant ecological properties
→ Identifiability thanks to shared random effects, but challenging validation
- Towards spatiotemporal, not purely spatial, modelling
→ Improve modelling of temporal dynamics
⚠ Requires disentangling complex observational/ecological dynamics
- Could we implement shared latent processes in other learning algorithms? (GAMs, ANNs, Random Forests...)


Discussion: Bias and uncertainty reduction



Citizen Science

- Checklist data, such as eBird, allow generating pseudo-absences, but many opportunistic datasets are less structured
- Data fusion of opportunistic and structured data in *Integrated Species Distribution Models* is crucial (Fithian et al 2015; Isaac et al 2020)
- Collecting additional exhaustive field data may be necessary
→ Explore optimal sampling design through simulation studies?

Discussion: Opportunities for ecological extreme-value analysis

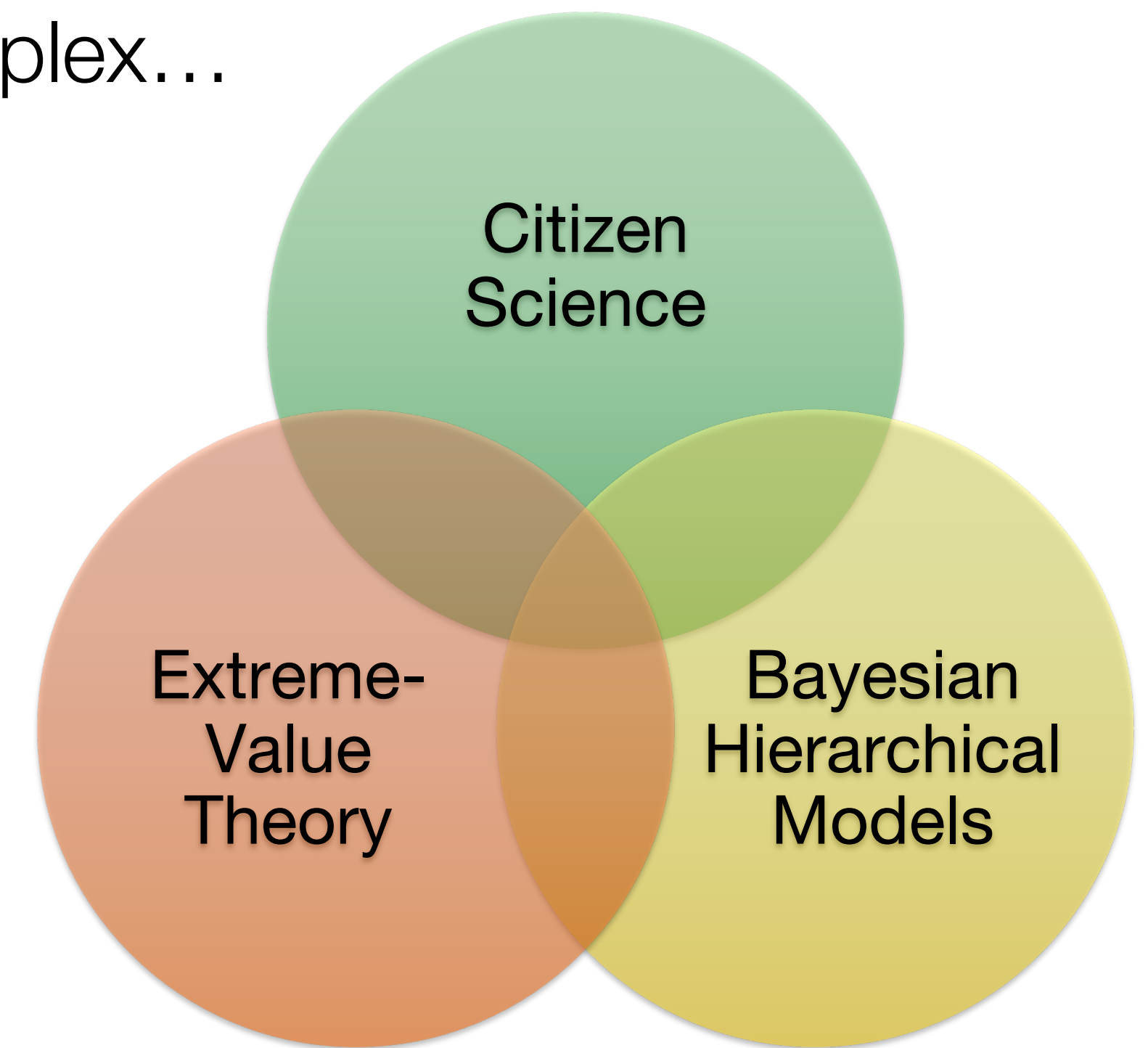


Extreme-Value Theory

- EVT generally less relevant for discrete data but promising for modelling extreme phenological events, such as first arrivals
- EVT widely used for extreme climate and environmental events
 - Such events can drive strong species population shifts
 - Focus on specific events, not only long-term climate averages
 - Probabilities and simulation for high-impact events

Outlook

- Rather basic handling of covariates and time trends in our model could be improved
- Extrapolated predictions could be validated using hold-out data by artificially reducing observational effort during training
- Ecological datasets: *Small Data* and *Big Data*, but always complex...
 - Wide opportunities for modelling and decision support
 - An exciting playground for statisticians!



Food for thought

This work:

Koh, Opitz (2024). Extreme-value modelling of migratory bird arrival dates: Insights from citizen science data. Journal of the Royal Statistical Society, Series A (Statistics in Society).

Other literature:

- Adjei et al. (2023). A structural model for the process of collecting biodiversity data. Authorea Preprints.
- Adjei et al. (2023). The Point Process Framework for Integrated Modelling of Biodiversity Data. arXiv:2311.06755.
- Belmont et al. (2024). Spatio-temporal Occupancy Models with INLA. arXiv:2403.10680.
- Coles (2001). An introduction to statistical modeling of extreme values. Springer.
- Diggle et al. (2010). Geostatistical inference under preferential sampling. Journal of the Royal Statistical Society Series C: Applied Statistics.
- Fithian et al. (2015). Bias correction in species distribution models: pooling survey and collection data for multiple species. Methods in Ecology and Evolution.
- Gelfand & Shirota (2019). Preferential sampling for presence/absence data and for fusion of presence/absence data with presence-only data. Ecological Monographs.
- Isaac et al. (2020). Data integration for large-scale models of species distributions. Trends in Ecology & Evolution.
- Lindgren et al. (2024). *inlabru*: software for fitting latent Gaussian models with non-linear predictors. arXiv:2407.00791.
- Tang et al. (2021). Modeling spatially biased citizen science effort through the eBird database. Environmental and Ecological Statistics.
- Wijeyakulasuriya et al. (2024). Modeling First Arrival of Migratory Birds Using a Hierarchical Max-Infinitely Divisible Process. Journal of Agricultural, Biological and Environmental Statistics.