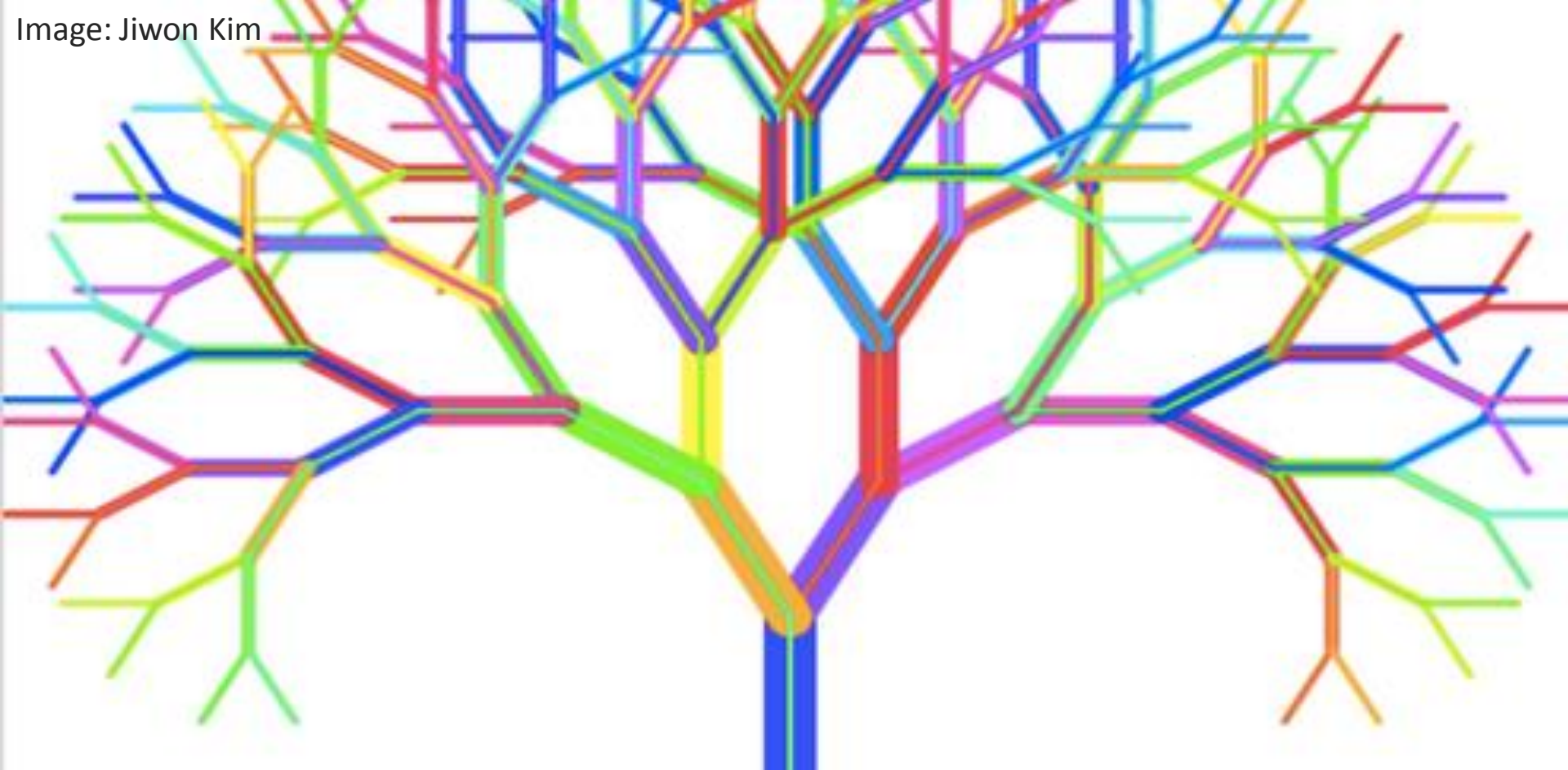


Image: Jiwon Kim



Can we use random forests for spatiotemporal CPUE modeling?

BRIAN STOCK, ERIC WARD, BRICE SEMMENS



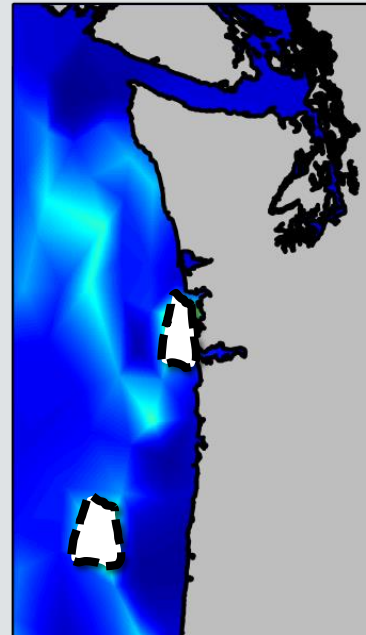
What we want (from Rick Methot)

- ★ Fast (coding vs. runtime vs. interpretation)
- ★ Replicable (method well-defined, get same answer)
- ★ Robust (insensitive to distributional assumptions, outliers)
- ★ Predictive ability (minimal errors, fill in space/time gaps)
- ★ Covariate effects (nonlinear, interactions)
- Uncertainty estimates (with known properties)
- ⊘ Specifiable structure (e.g. correlation through time, biology)
- ⊘ Unbiased (relative vs. absolute abundance)

What we want (from Rick Methot)

- ★ Fast
- ★ Replicable
- ★ Robust
- ★ Predictive ability
- ★ Covariate effects
- Uncertainty estimates
- ⊘ Specifiable structure
- ⊘ Unbiased

Story 1: Bycatch hotspots

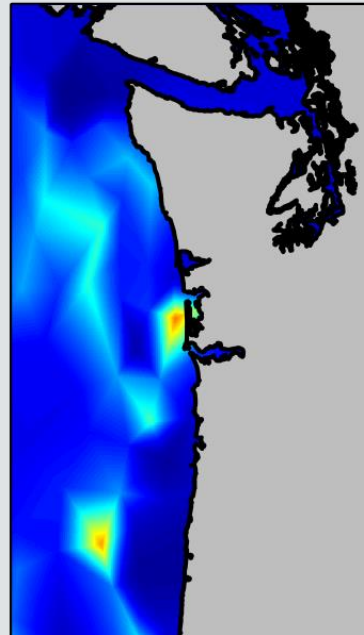


What we want (from Rick Methot)

- ★ Fast
- ★ Replicable
- ★ Robust
- ★ Predictive ability
- ★ Covariate effects
- Uncertainty estimates
- ⊘ Specifiable structure
- ⊘ Unbiased

Story 2: Total bycatch estimation

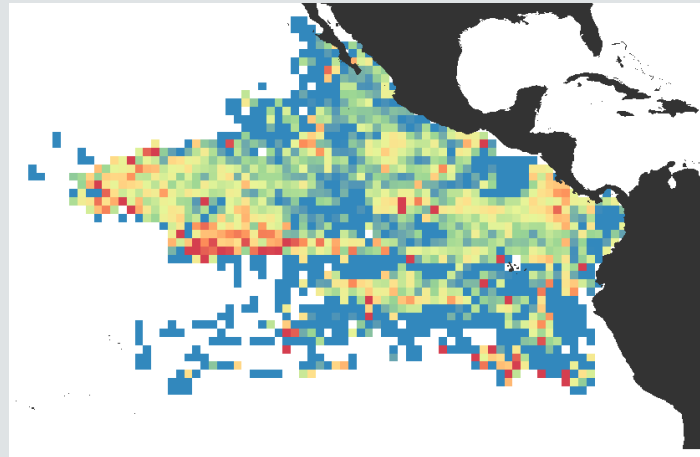
Σ



What we want (from Rick Methot)

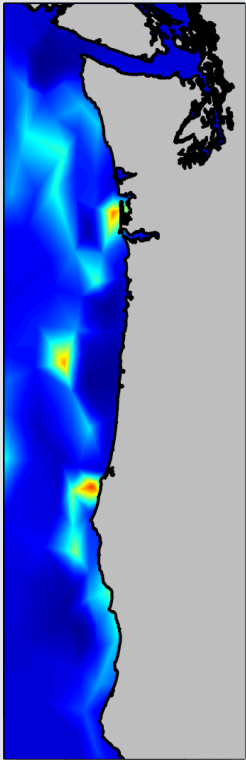
- ★ Fast
- ★ Replicable
- ★ Robust
- ★ Predictive ability
- ★ Covariate effects
- Uncertainty estimates
- ⊘ Specifiable structure
- ⊘ Unbiased

Story 3: CPUE standardization



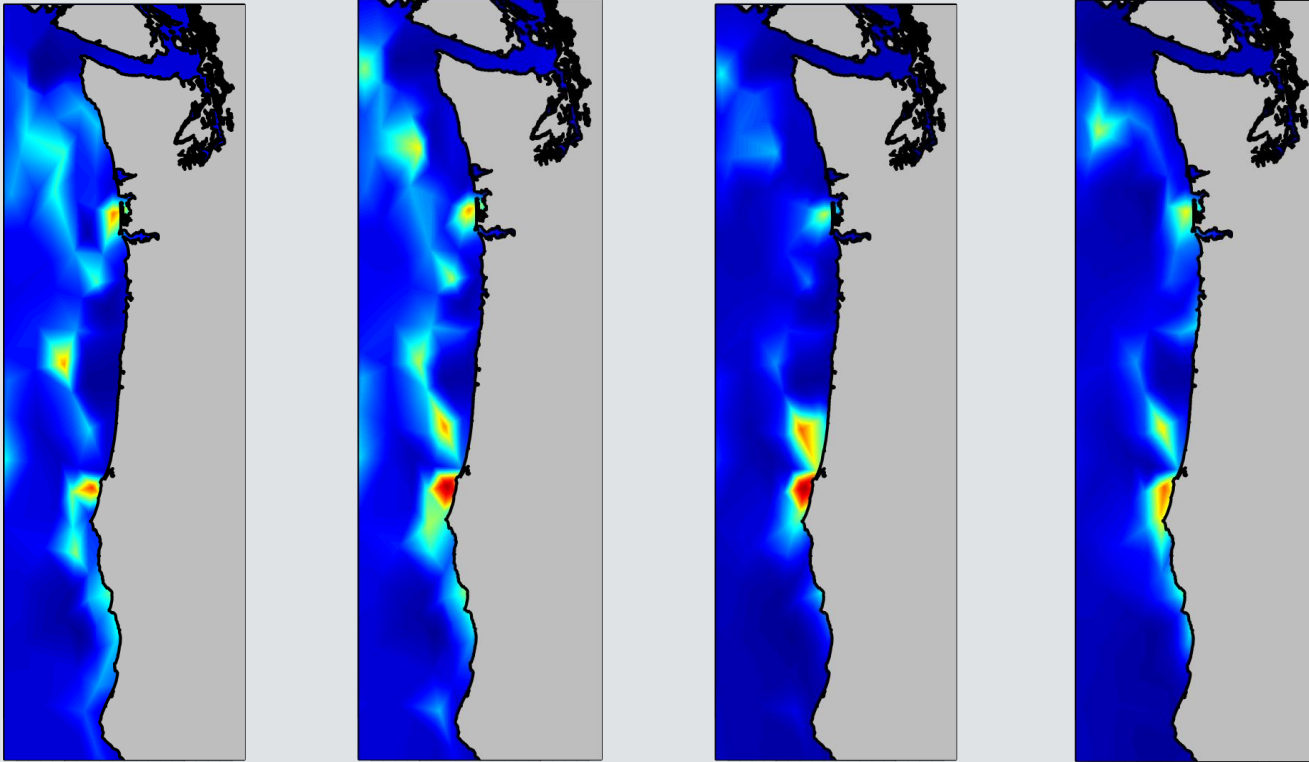
Tools for dynamic management

Need map of bycatch “risk”



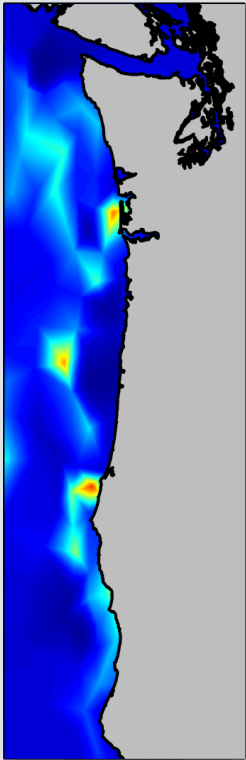
Tools for dynamic management

Need map of bycatch “risk”



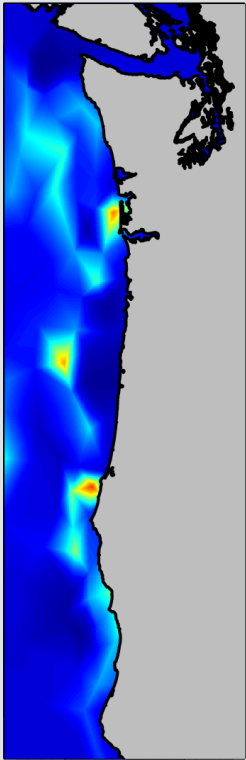
Tools for dynamic management

Need map of bycatch “risk”



- temperature
- depth
- substrate
- spatial field

Q: Which spatial model is best?



- temperature
- depth
- substrate
- **spatial field**

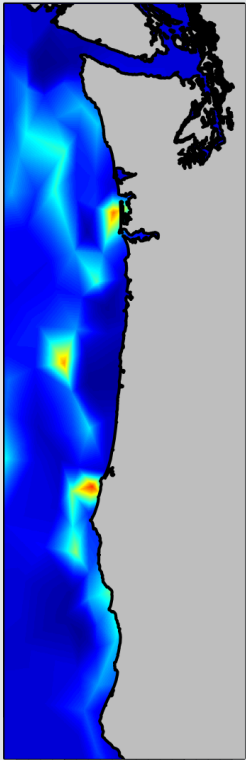
GLM

GAM

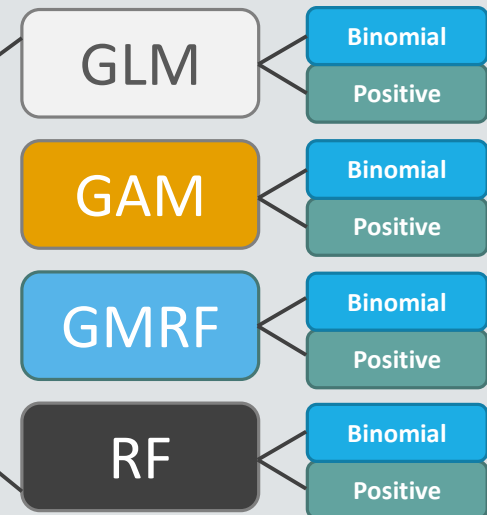
GMRF

RF

Q: Which spatial model is best?



- temperature
- depth
- substrate
- **spatial field**



What are these models exactly?

GLM

obs \sim environmental predictors (temp, depth, ...)

$$Y_i \sim \text{Binomial}(\text{logit}^{-1}[\mathbf{X}_i\boldsymbol{\beta}])$$

Binomial

$$Y_i \sim \text{Gamma}(e^{\mathbf{X}_i\boldsymbol{\beta}}, \nu)$$

Positive

GAM

GMRF

RF

What are these models exactly?

GLM

obs ~ environmental predictors (temp, depth, ...)

GAM

obs ~ environmental predictors + s(lat,lon)

GMRF

RF

What are these models exactly?

GLM

obs \sim environmental predictors (temp, depth, ...)

GAM

obs \sim environmental predictors + $s(\text{lat}, \text{lon})$

GMRF

obs \sim environmental predictors + $MVN(0, \Sigma)$

RF

What are these models exactly?

GLM

obs \sim environmental predictors (temp, depth, ...)

GAM

obs \sim environmental predictors + $s(\text{lat}, \text{lon})$

GMRF

obs \sim environmental predictors + $MVN(0, \Sigma)$

RF

obs \sim environmental predictors + lat + lon

Fisheries observer data

U.S. West Coast
Groundfish
Trawl



Hawaii
Swordfish
Longline



Generally:

GLM

<

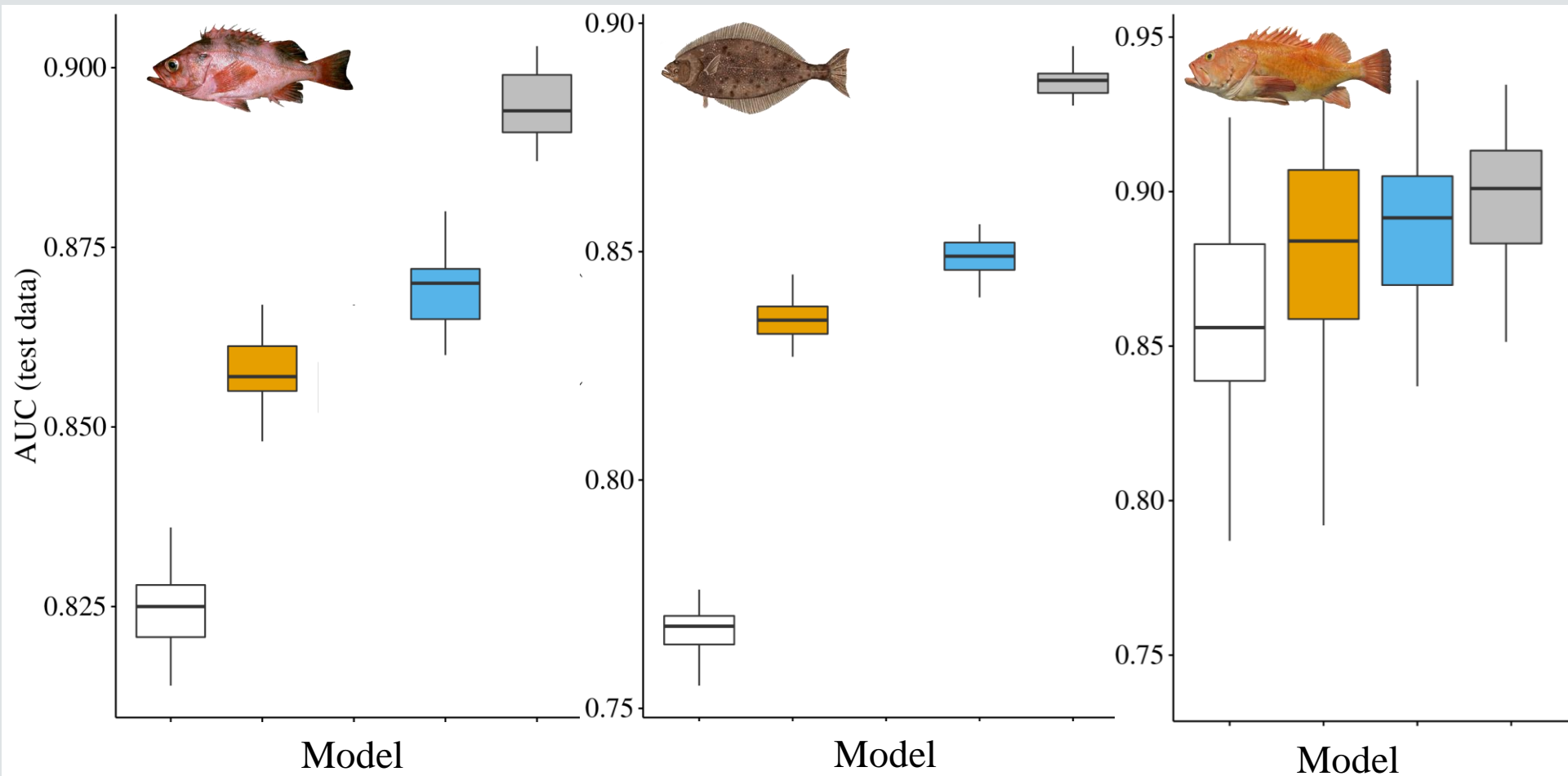
GAM

<

GMRF

<

RF



Generally:

GLM

<

GAM

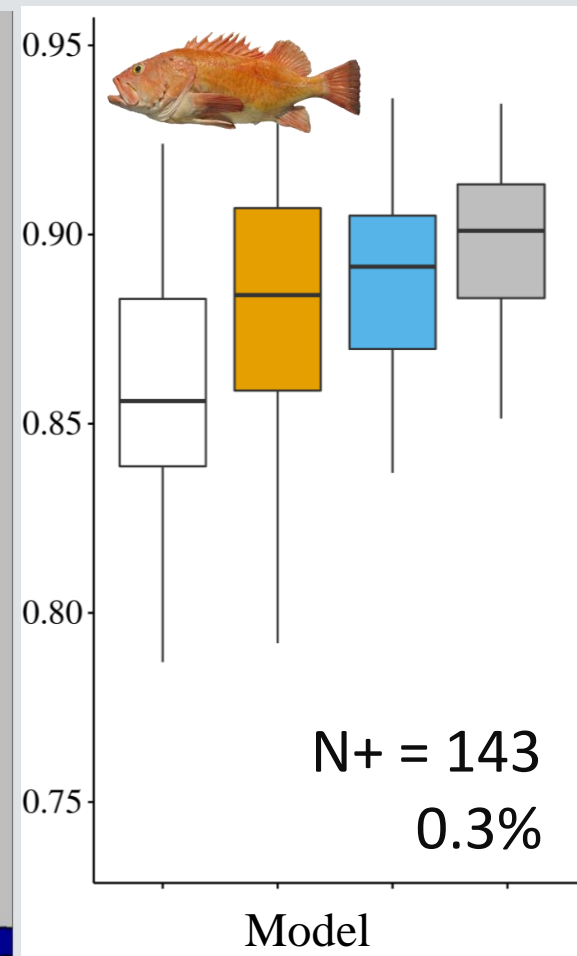
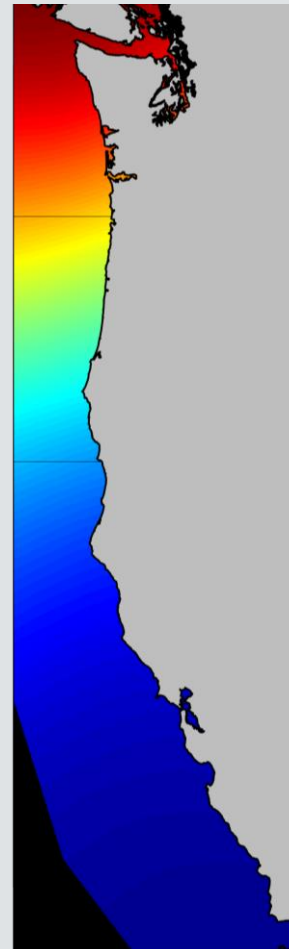
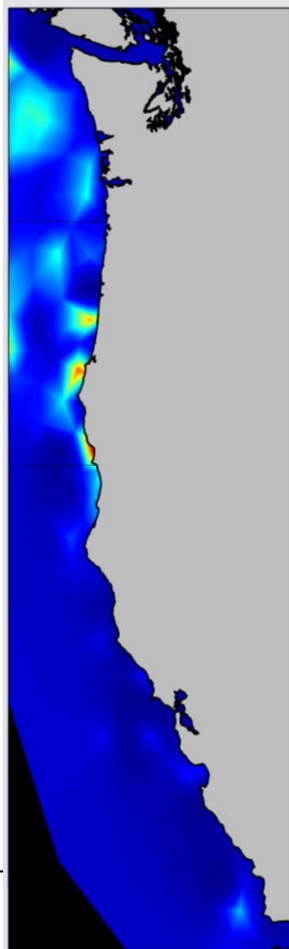
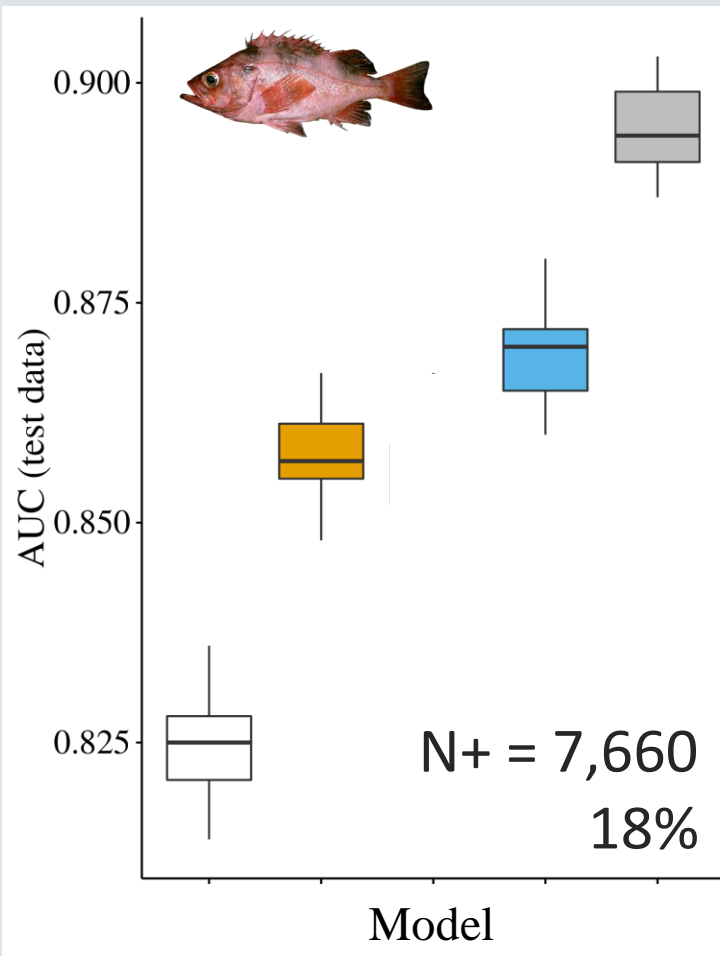
<

GMRF

<

RF

Less clear for rarer species



Generally:

GLM

<

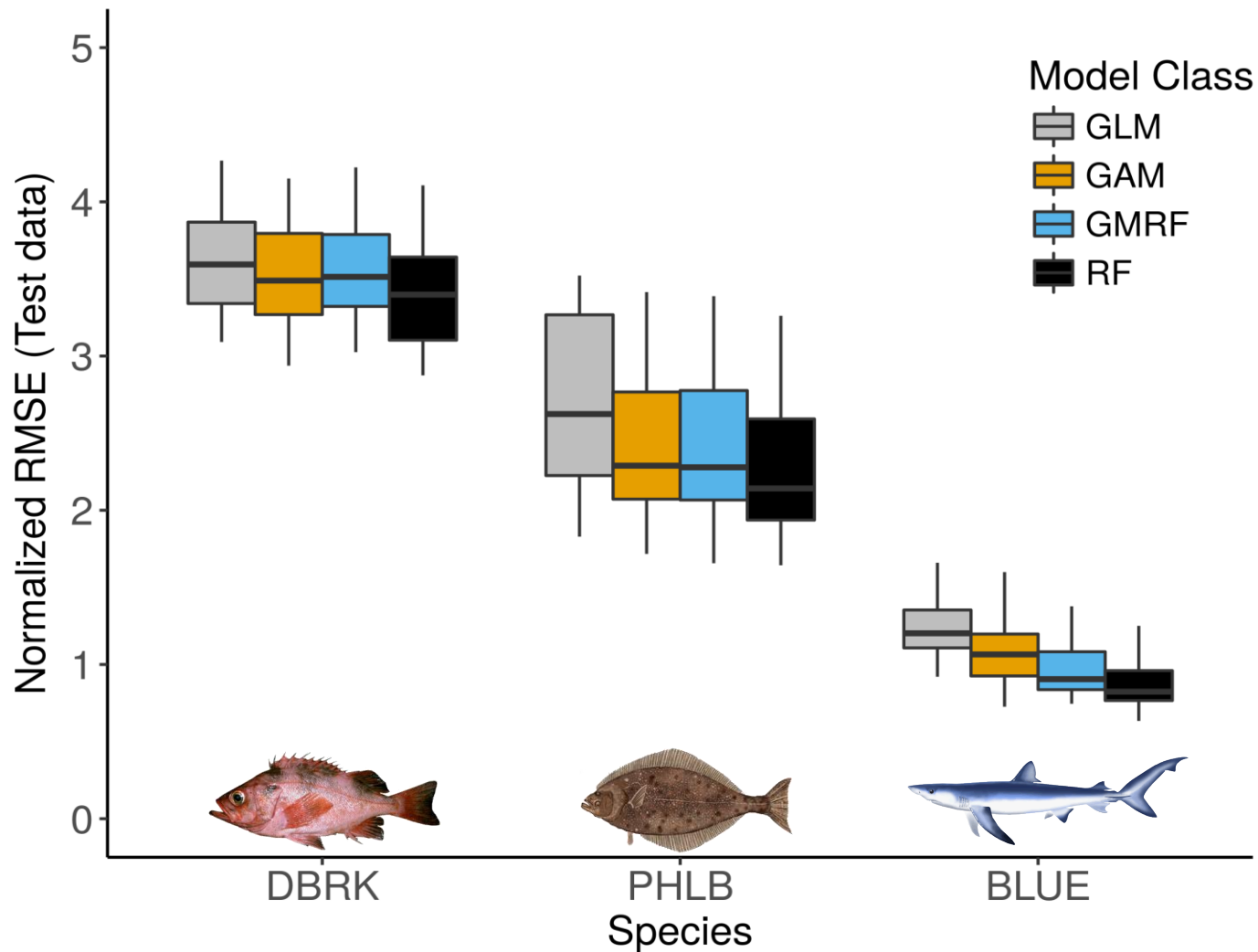
GAM

<

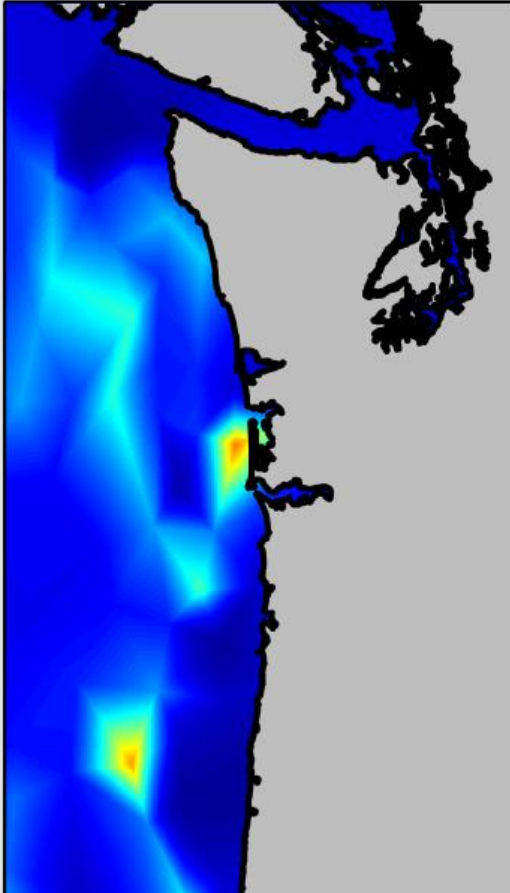
GMRF

<

RF



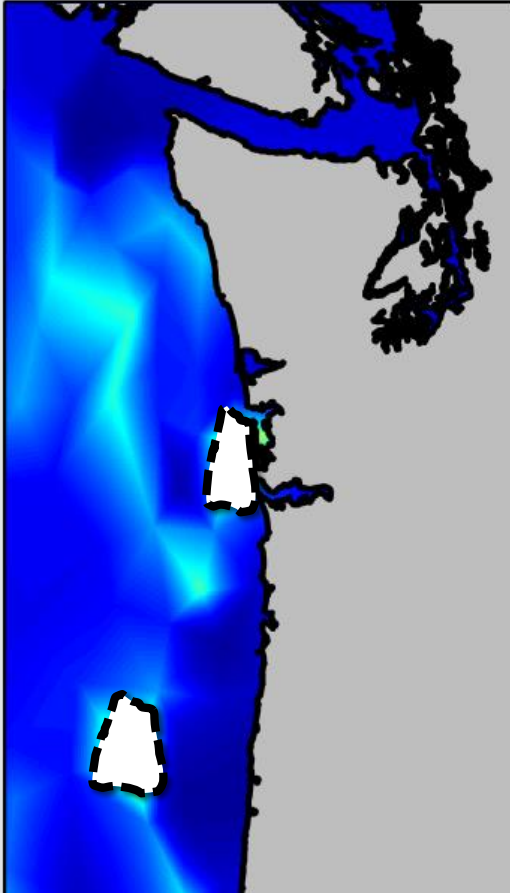
Q: How much bycatch can they prevent?



Crude management simulation:

1. Predict bycatch risk at test locations

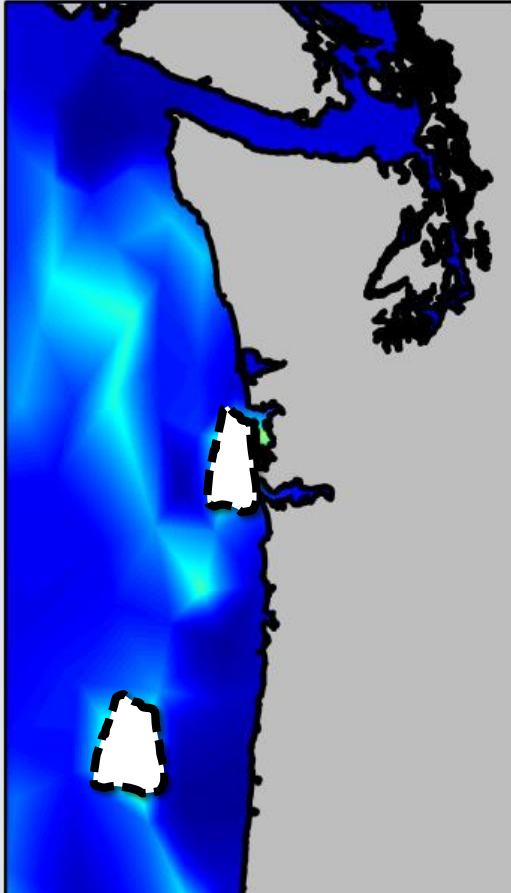
Q: How much bycatch can they prevent?



Crude management simulation:

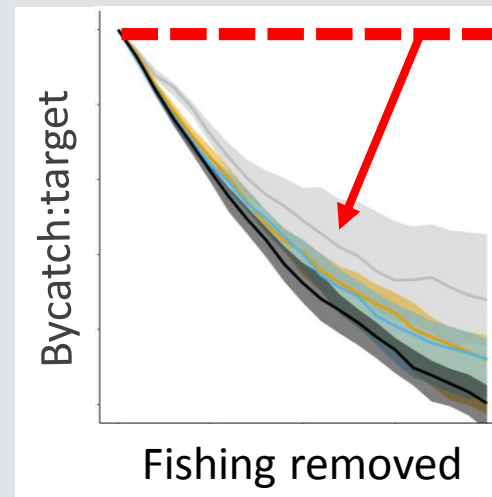
1. Predict bycatch risk at test locations
2. Remove X% of fishing effort with highest bycatch risk

Q: How much bycatch can they prevent?

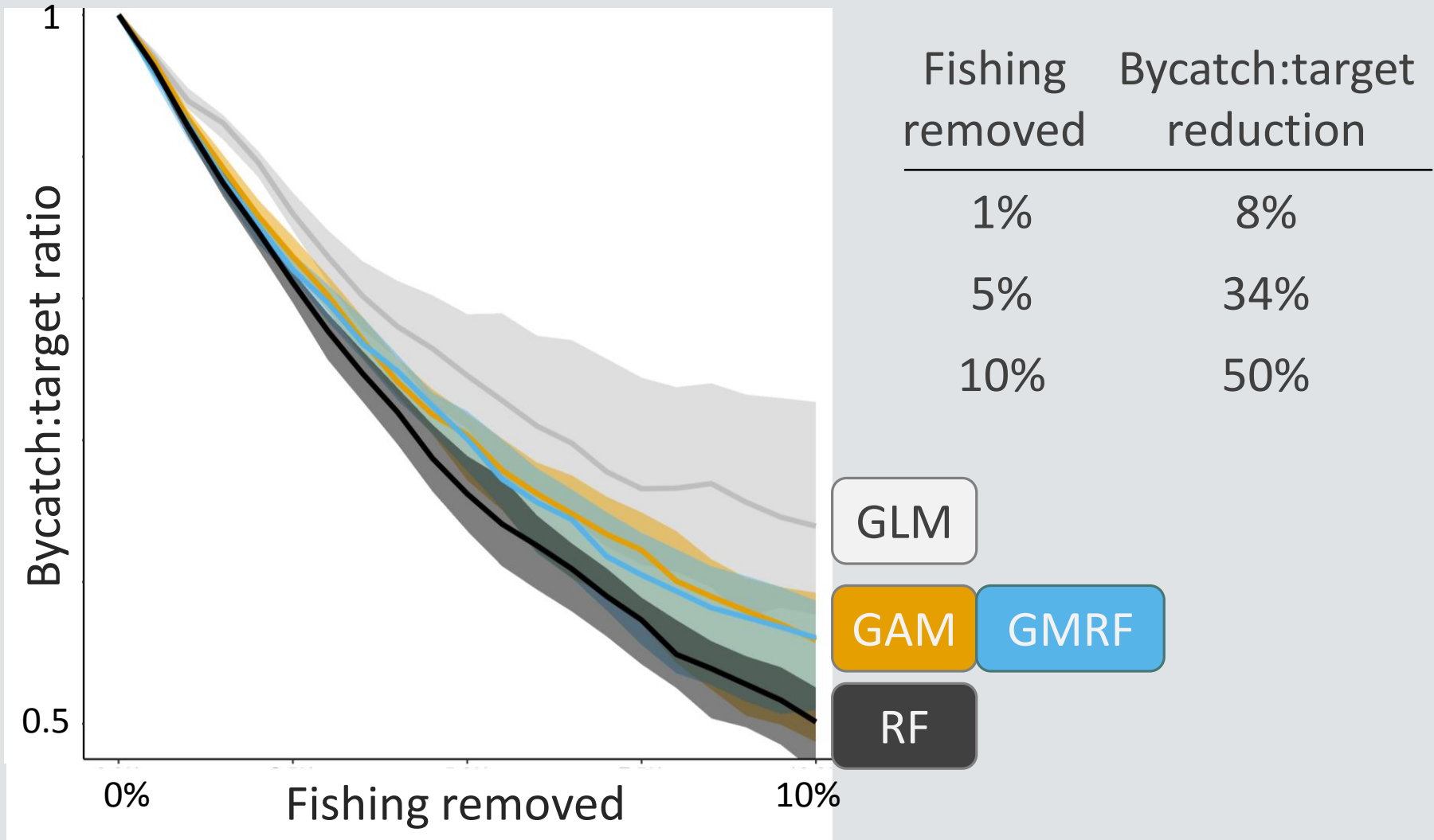


Crude management simulation:

1. Predict bycatch risk at test locations
2. Remove X% of fishing effort with highest bycatch risk
3. Calculate “prevented” bycatch and target catch (bycatch:target ratio)



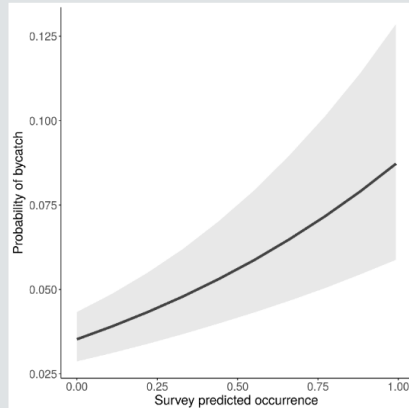
Q: How much bycatch can they prevent?



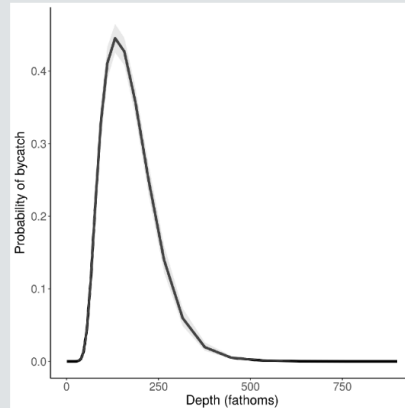
Covariate effects



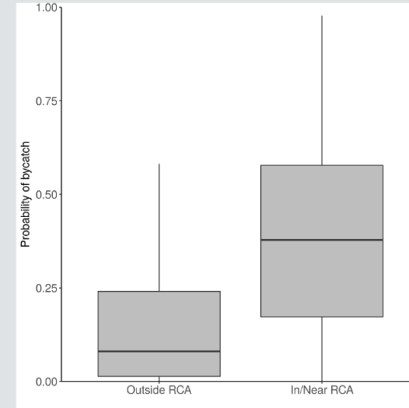
GMRF



PredOccSurvey



Depth

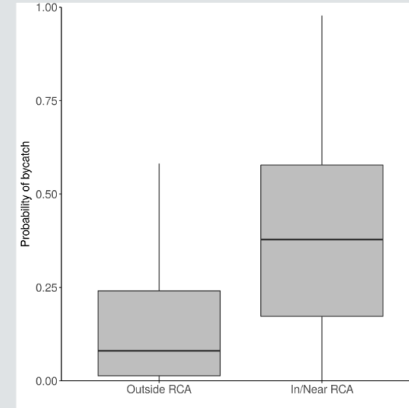
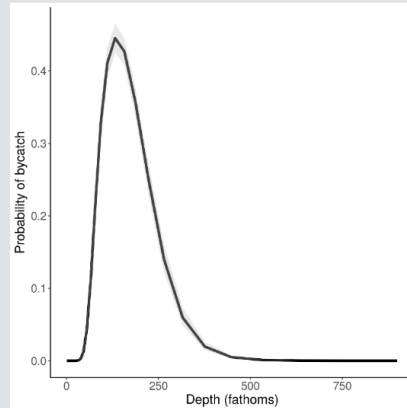
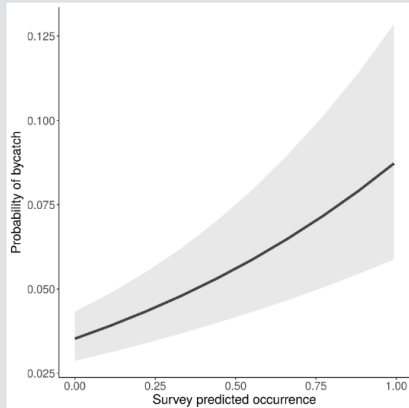


In/near RCA

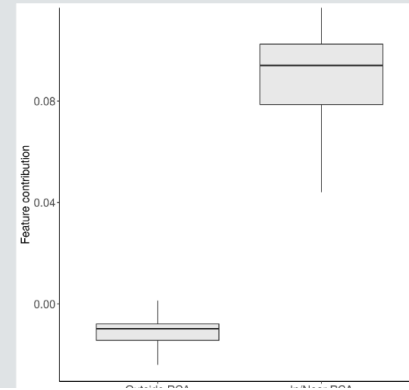
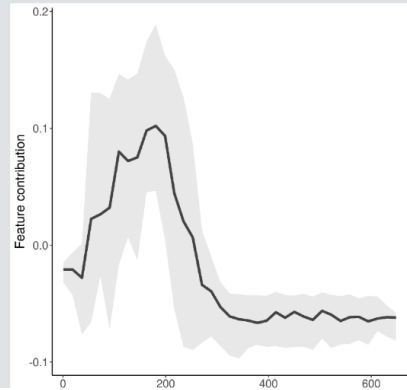
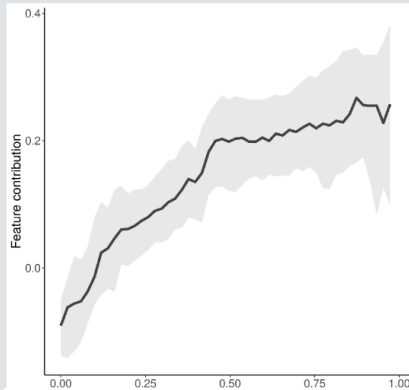
Covariate effects



GMRF



RF

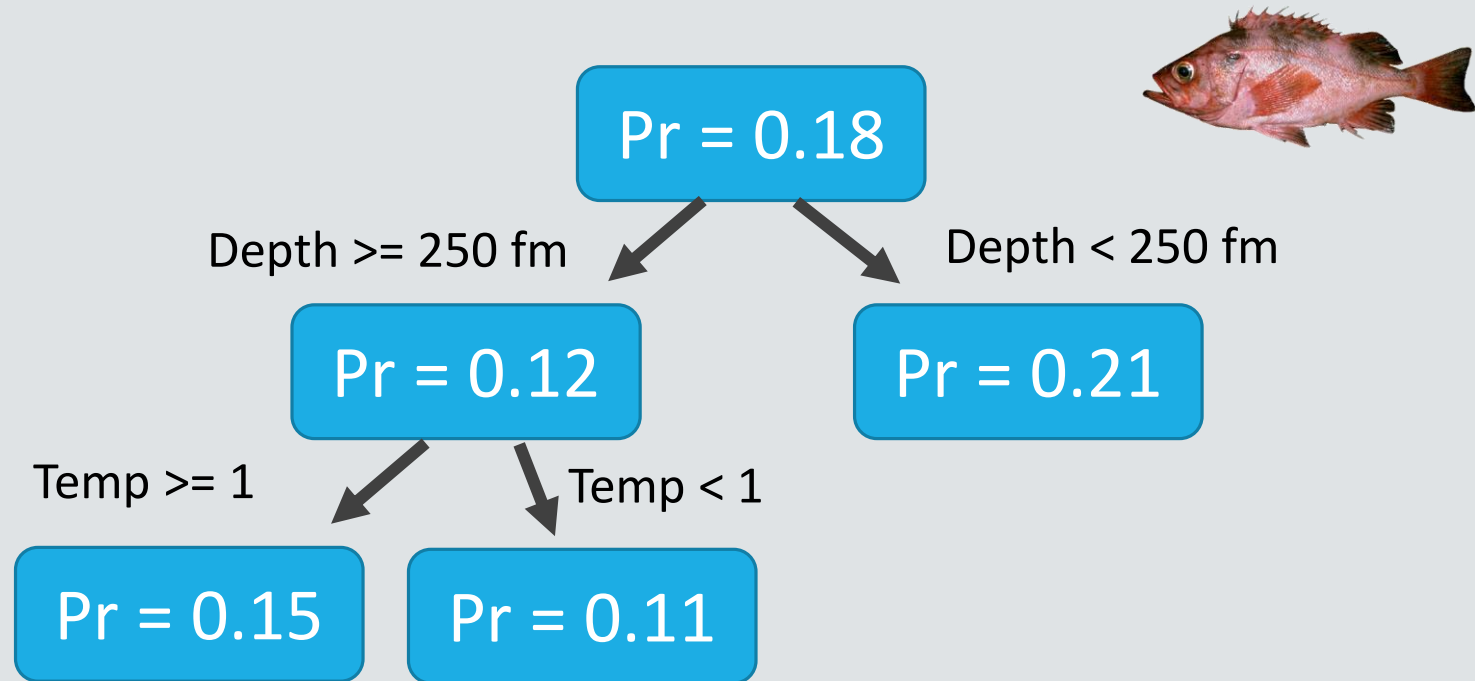


How?

How do random forests work?

Single decision tree:

Low bias, high variance model (overfit)

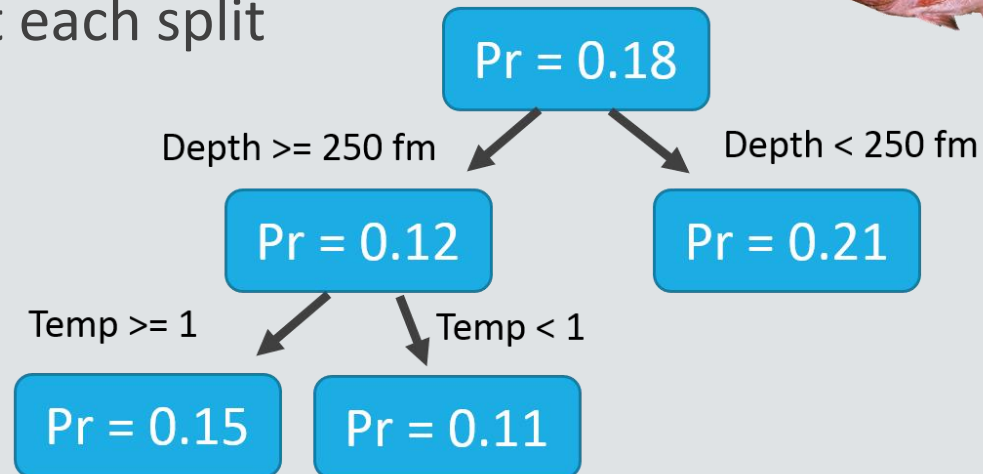


How do random forests work?

Idea: average across many, uncorrelated trees

$$E[MSE] = Model\ Bias^2 + Model\ Variance + noise$$

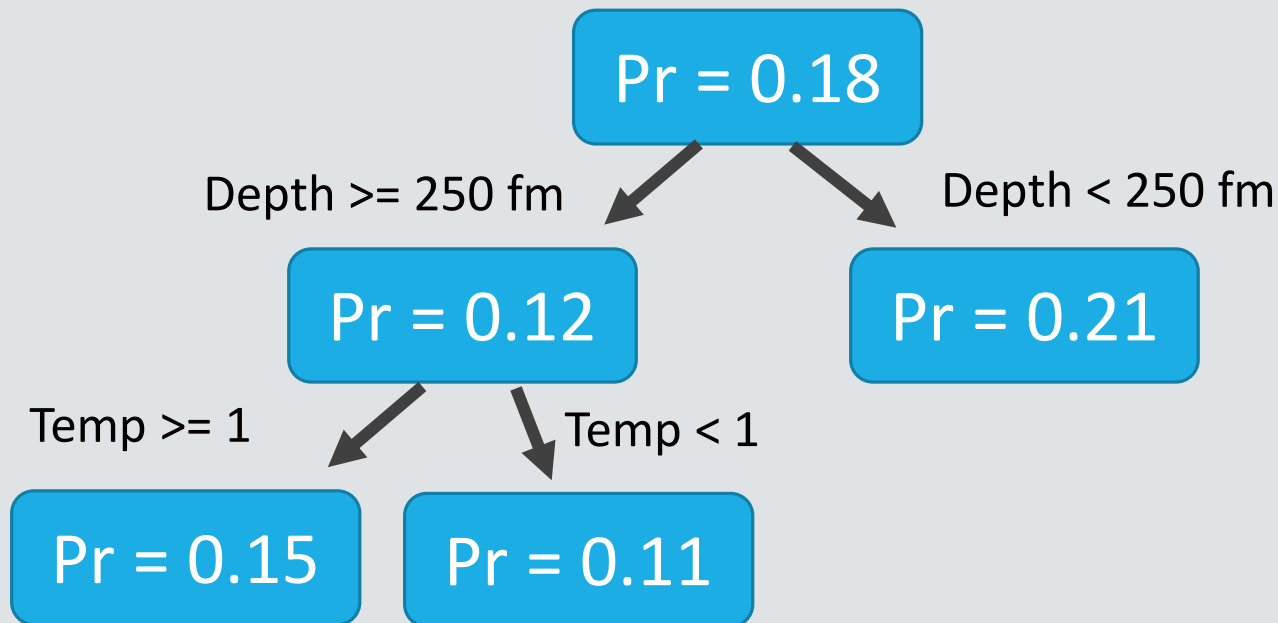
1. *Bagging*: fit each tree on a **B**ootstrap sample (~63%) of the data, then **A**ggregate across trees (~1000+)
2. Only consider a *random subset* (~P/3) of *covariates* at each split



Covariate effects with RF



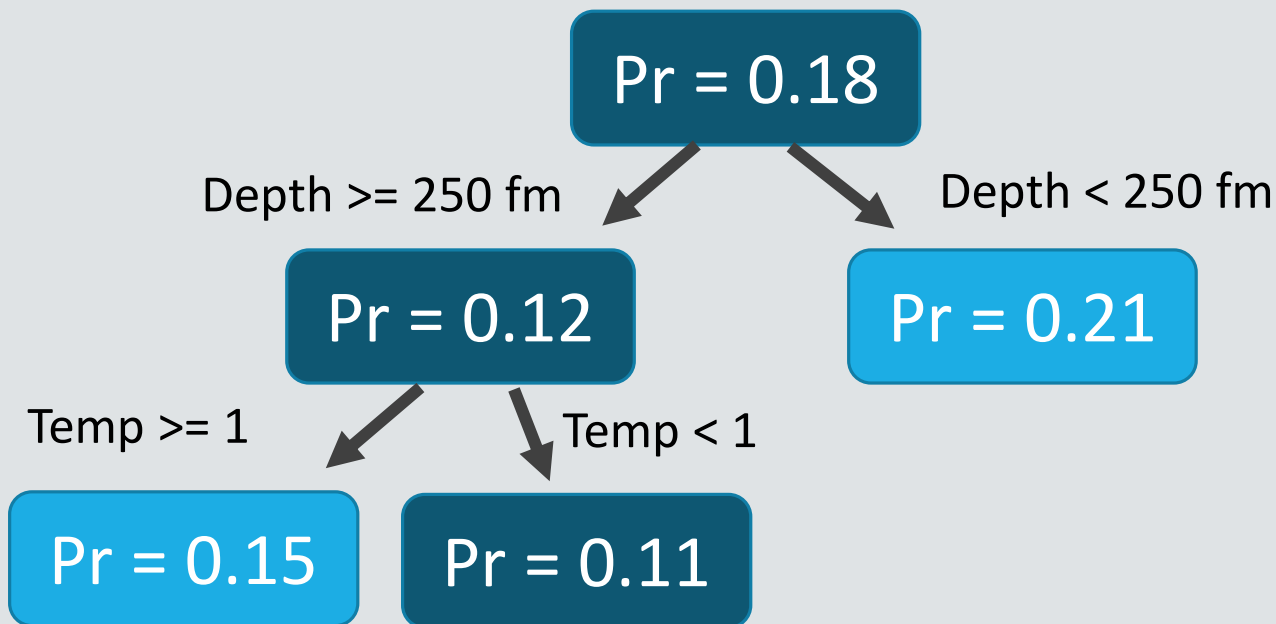
What is a “feature contribution”??



Covariate effects with RF



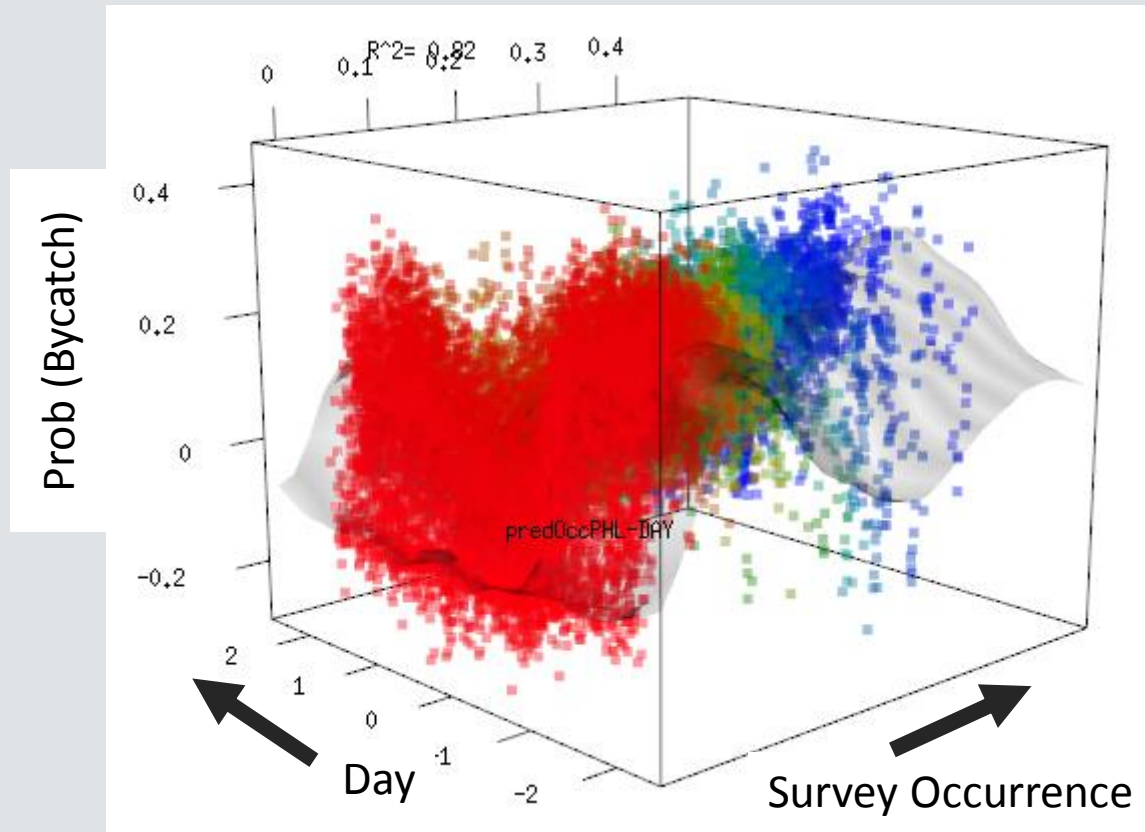
What is a “feature contribution”??



$$\text{Prediction}_i = 0.11 = 0.18 - 0.06 \text{ (Depth)} - 0.01 \text{ (Temp)}$$

Covariate interactions with RF

Catchability varies by Julian Day



1. Discussion

#2: Total bycatch estimates

Need estimates of total bycatch / discards

- Rarely observe 100% of fishing
- Often observe ~20%

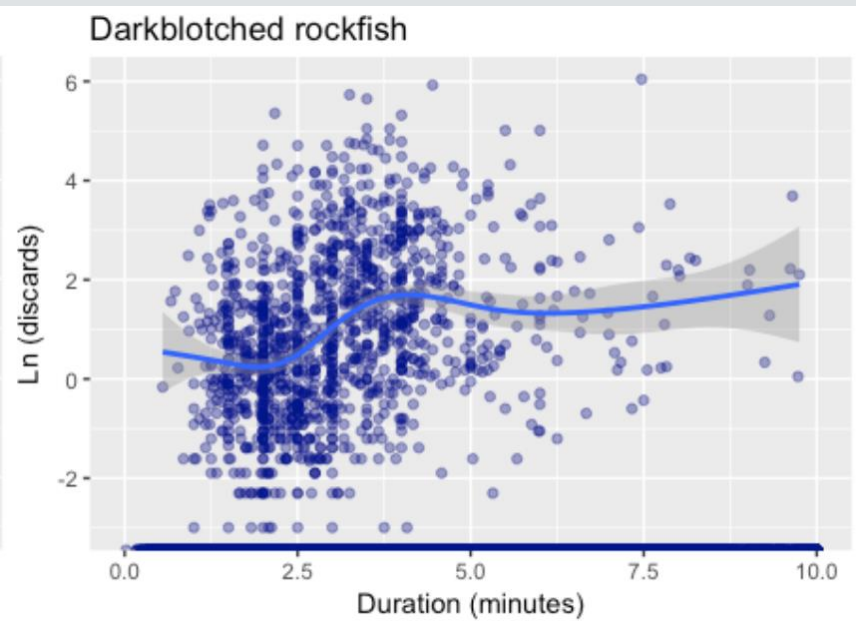
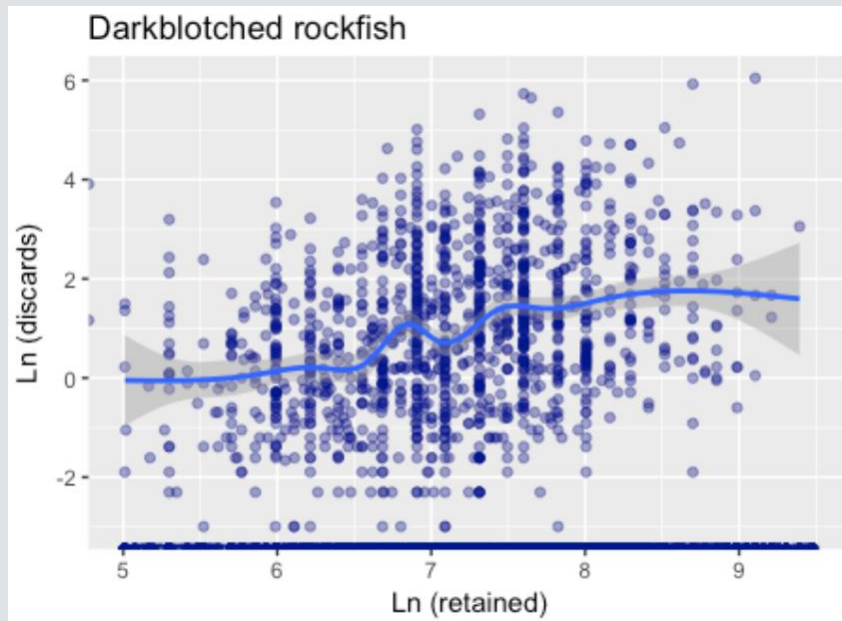
#2: Total bycatch estimates

“Ratio estimator”:

$$B_{unobs} = T_{unobs} \frac{B_{obs}}{T_{obs}}$$



Assumes bycatch prop. to target catch / effort

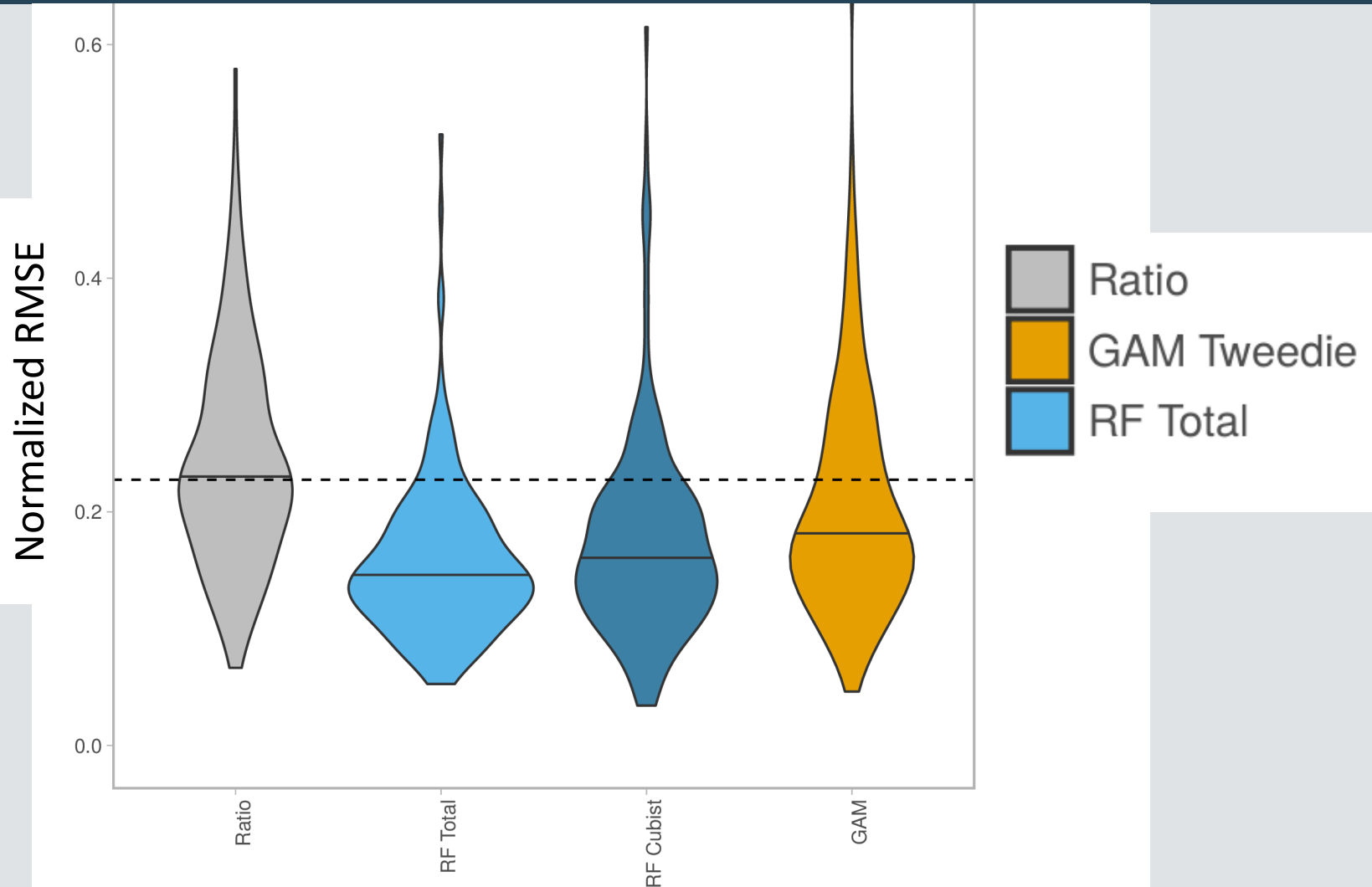


Use a spatial model instead

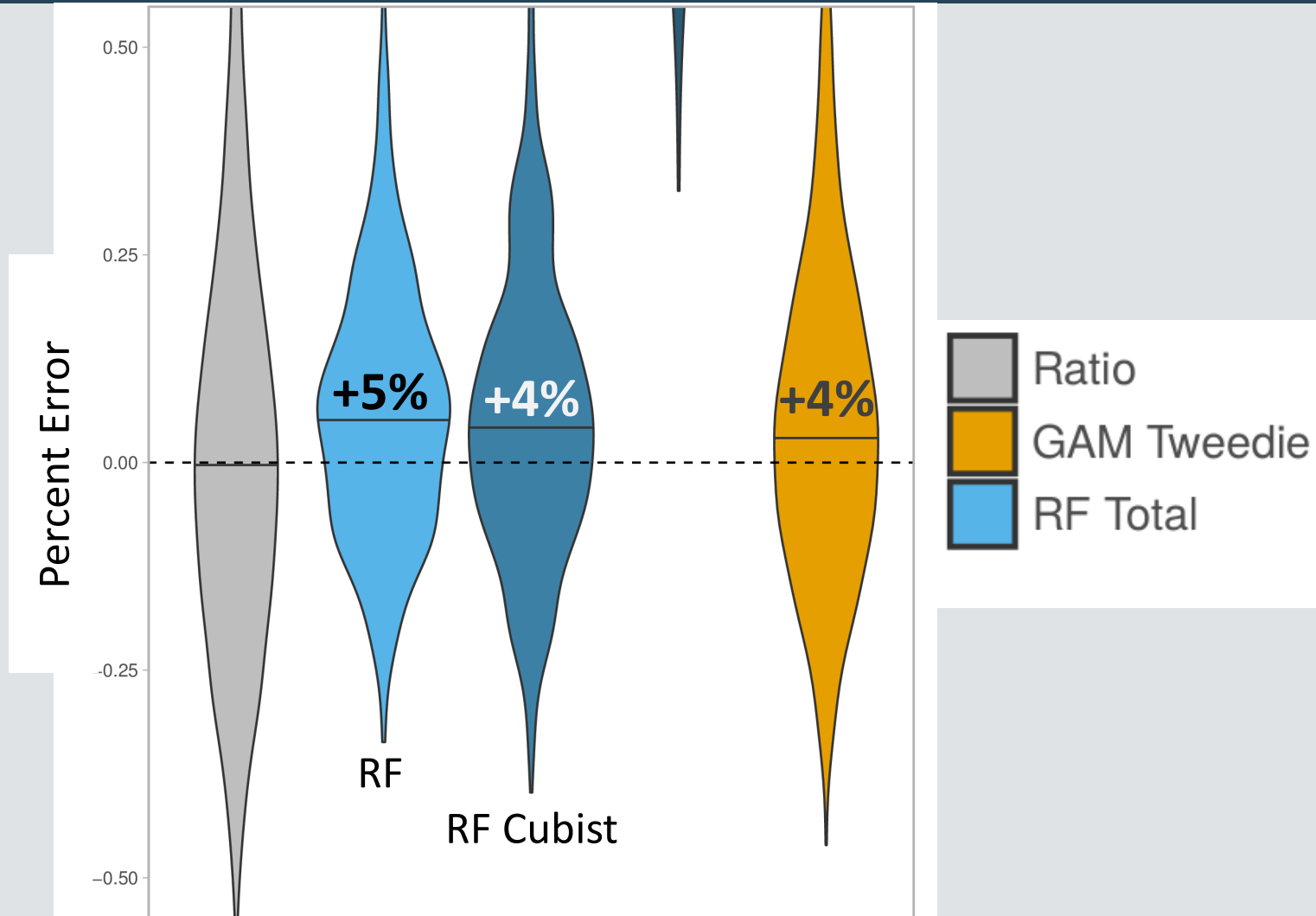
Cross-validation using dataset with 100% coverage:

1. Choose 20% observed trips
2. Fit spatial model
3. Predict at 80% unobserved
4. Compare $\text{sum}(\text{predictions})$ to ratio estimator

Spatial models = lower error

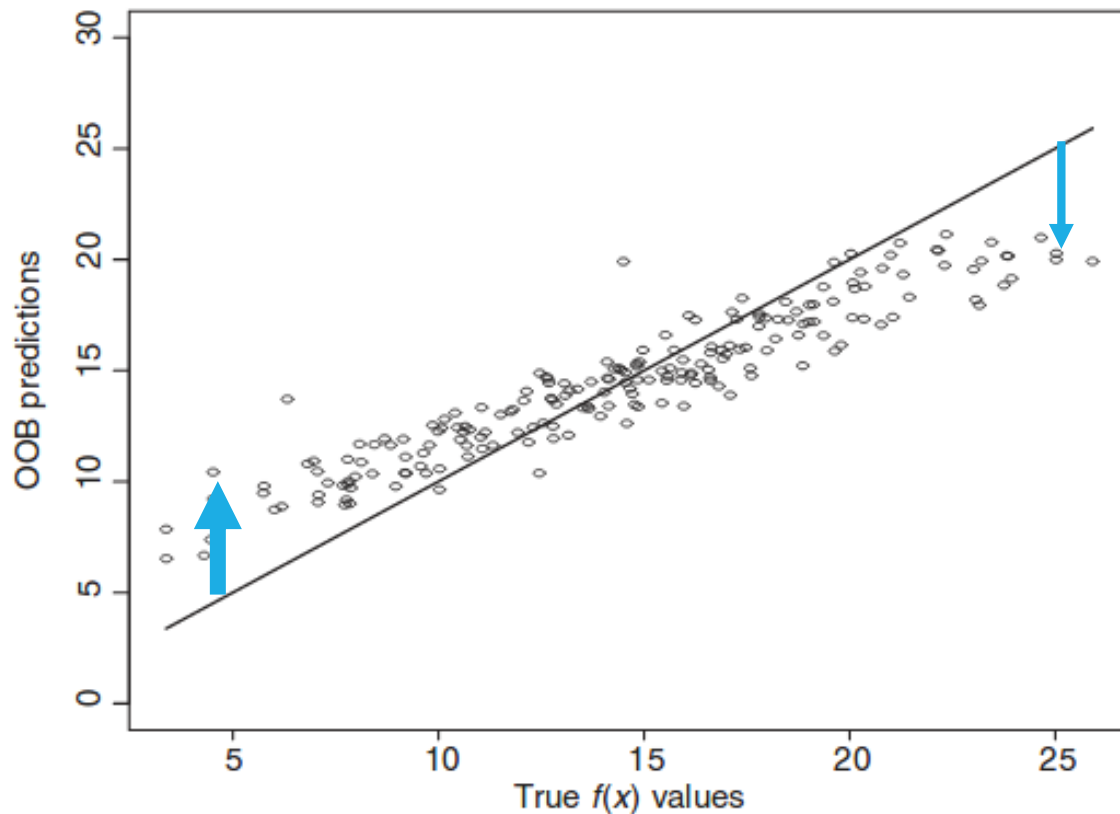


... bias in spatial model estimates



Why are random forests biased?

1. Extreme values modeled using average of less-extreme points →
Regression to the mean



2. Bycatch distribution
is right-skewed

Thoughts on RF bias

Bias correction methods:

- Fit linear model in nodes instead of mean ('Cubist')
- Fit second model on RF residuals (Xu 2013)



Bycatch estimates (*absolute* abundance) vs.
CPUE standardization (*relative* abundance)

#3: CPUE data

Eastern Pacific Ocean yellowfin tuna

- 2000-2009 catch + effort
- 1-deg lat/lon bins

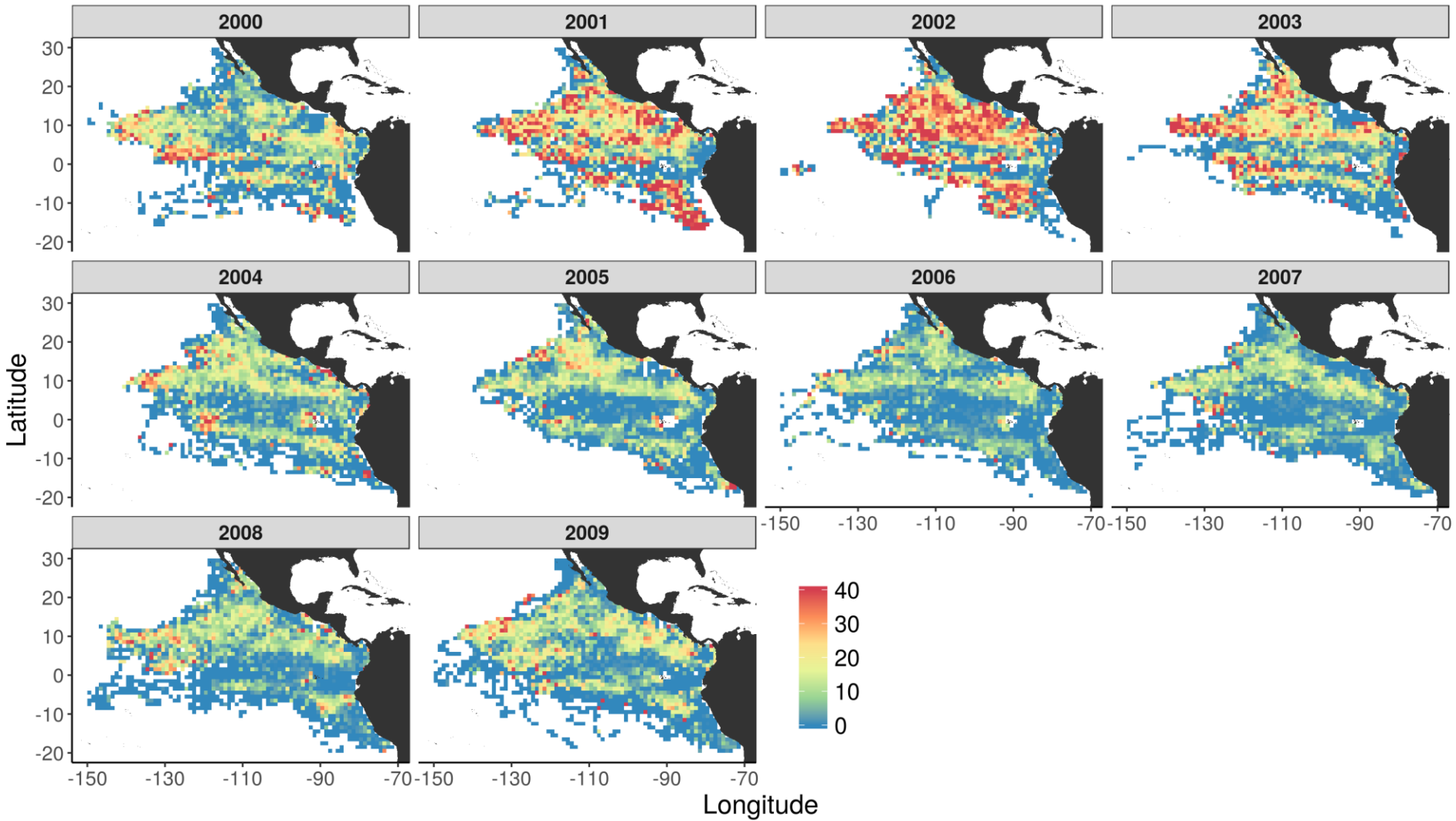
Model:

- 2000-2009 catch + effort
- 1-deg lat/lon bins

'ranger' `ranger(cpue ~ lat + lon + year, ...)`

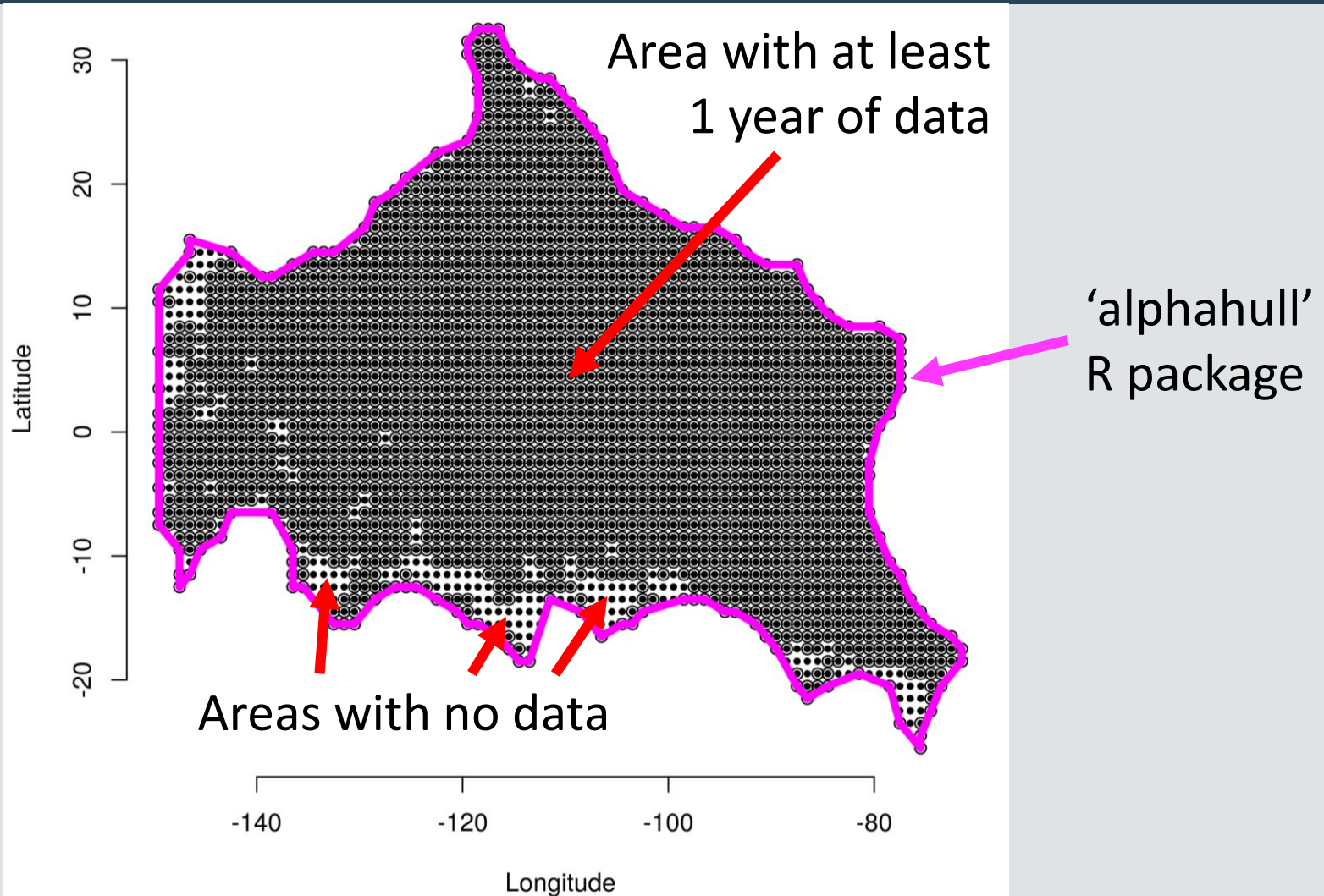
'grf' `regression_forest(dat[,covar], Y=dat$cpue, ...)`

CPUE data

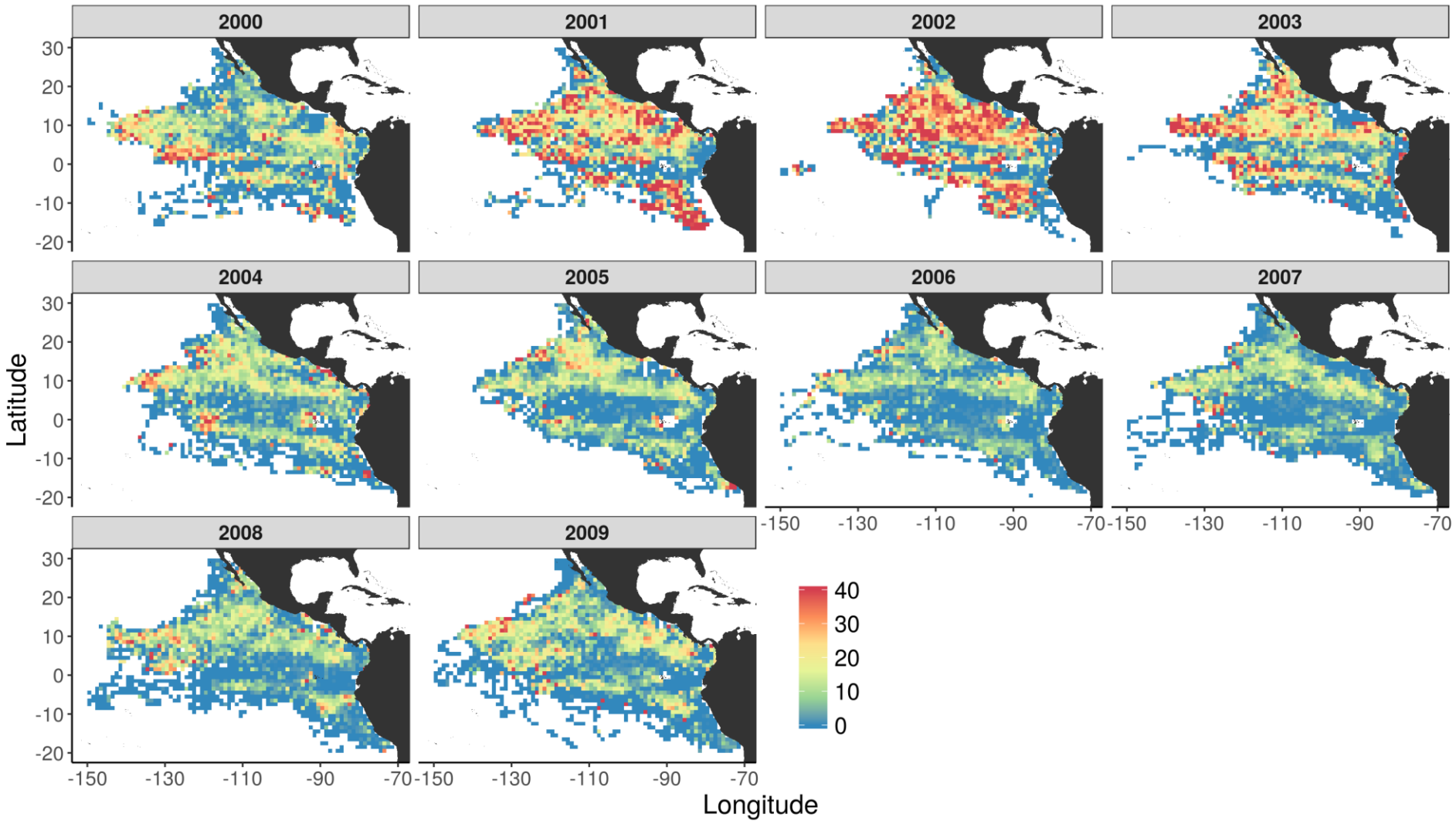


3. Methods

Create prediction grid

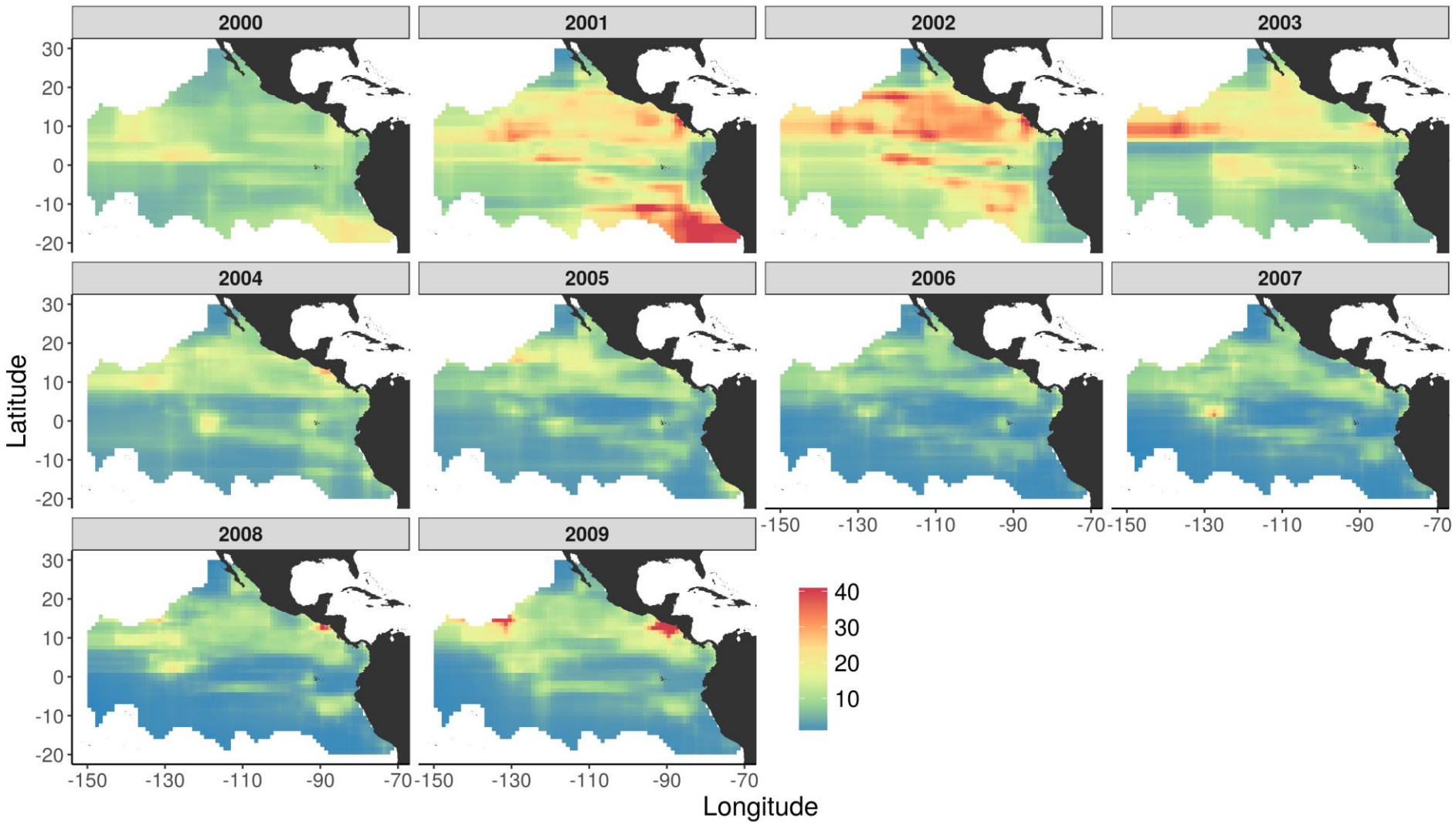


CPUE data



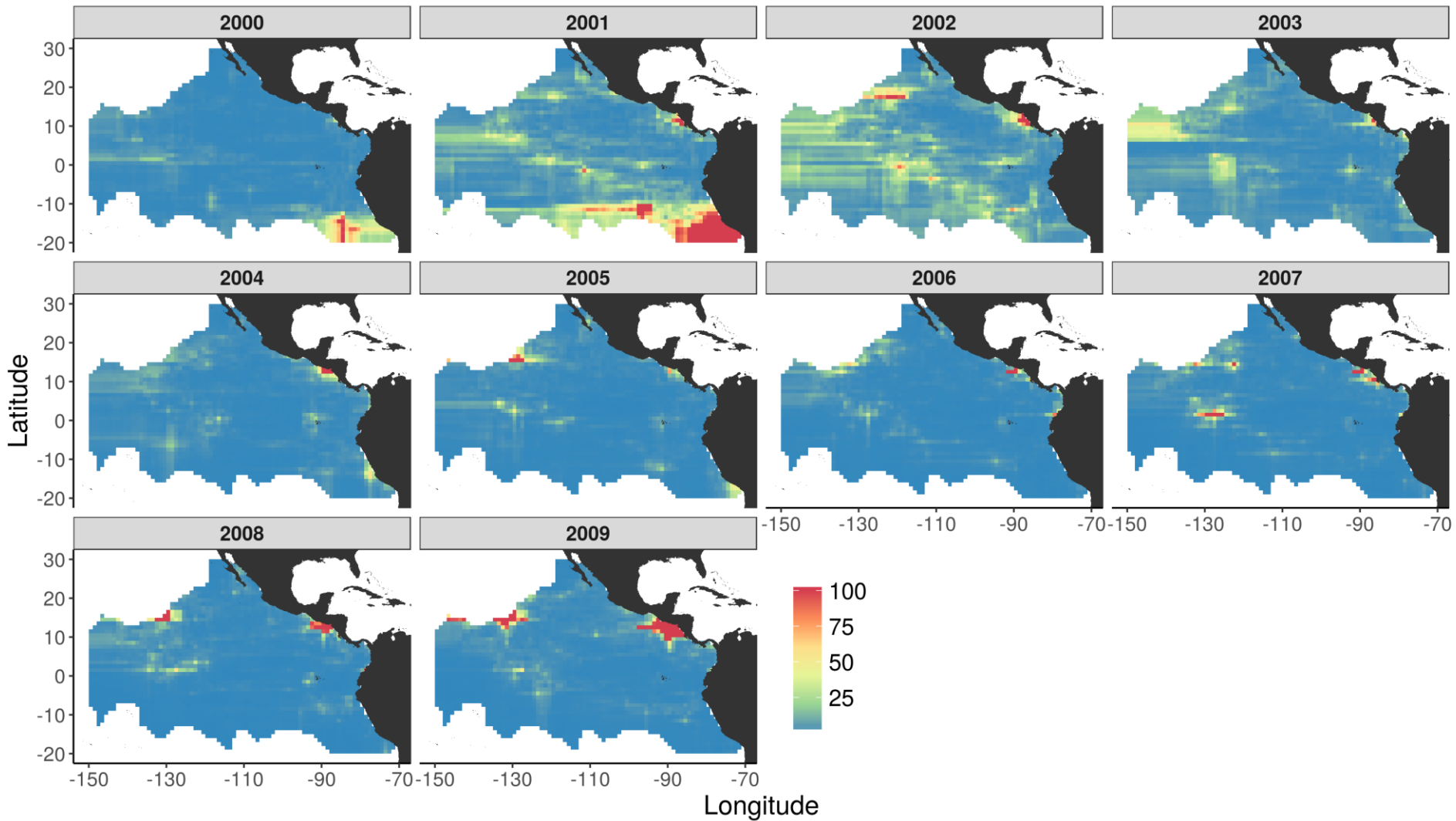
3. Methods

Predicted mean(CPUE)



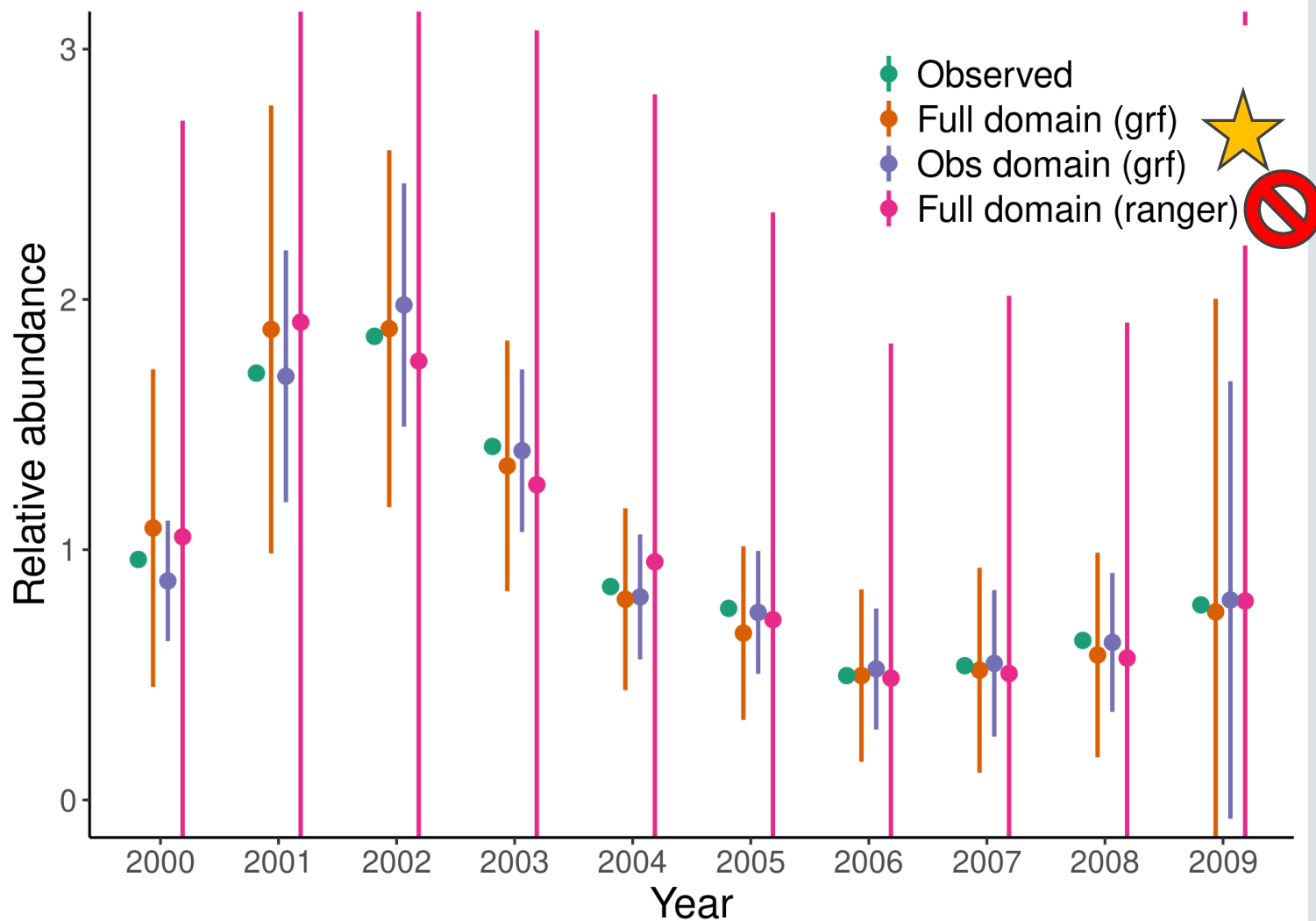
3. Results

Predicted Var(CPUE)

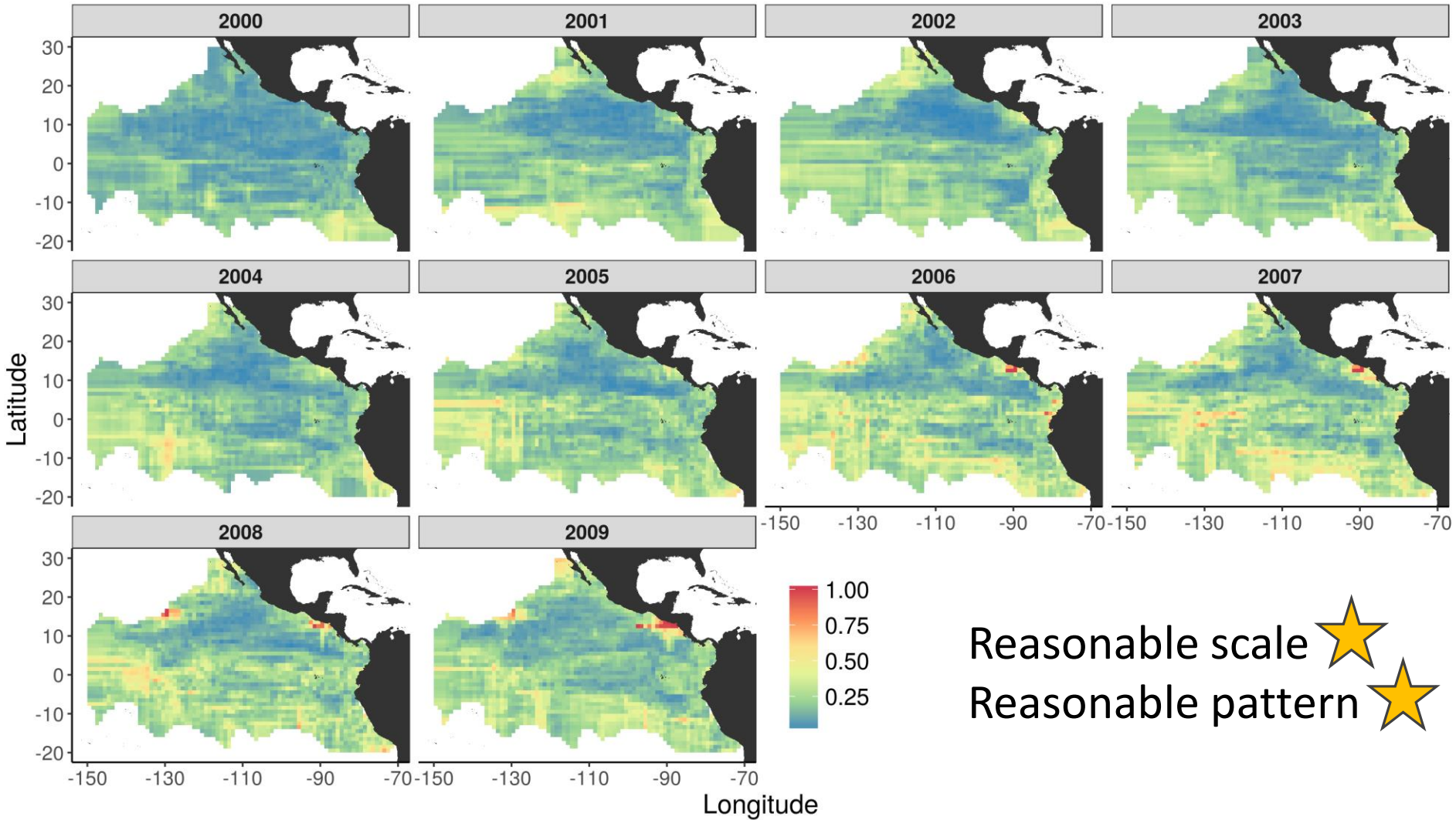


3. Results

Relative abundance trend

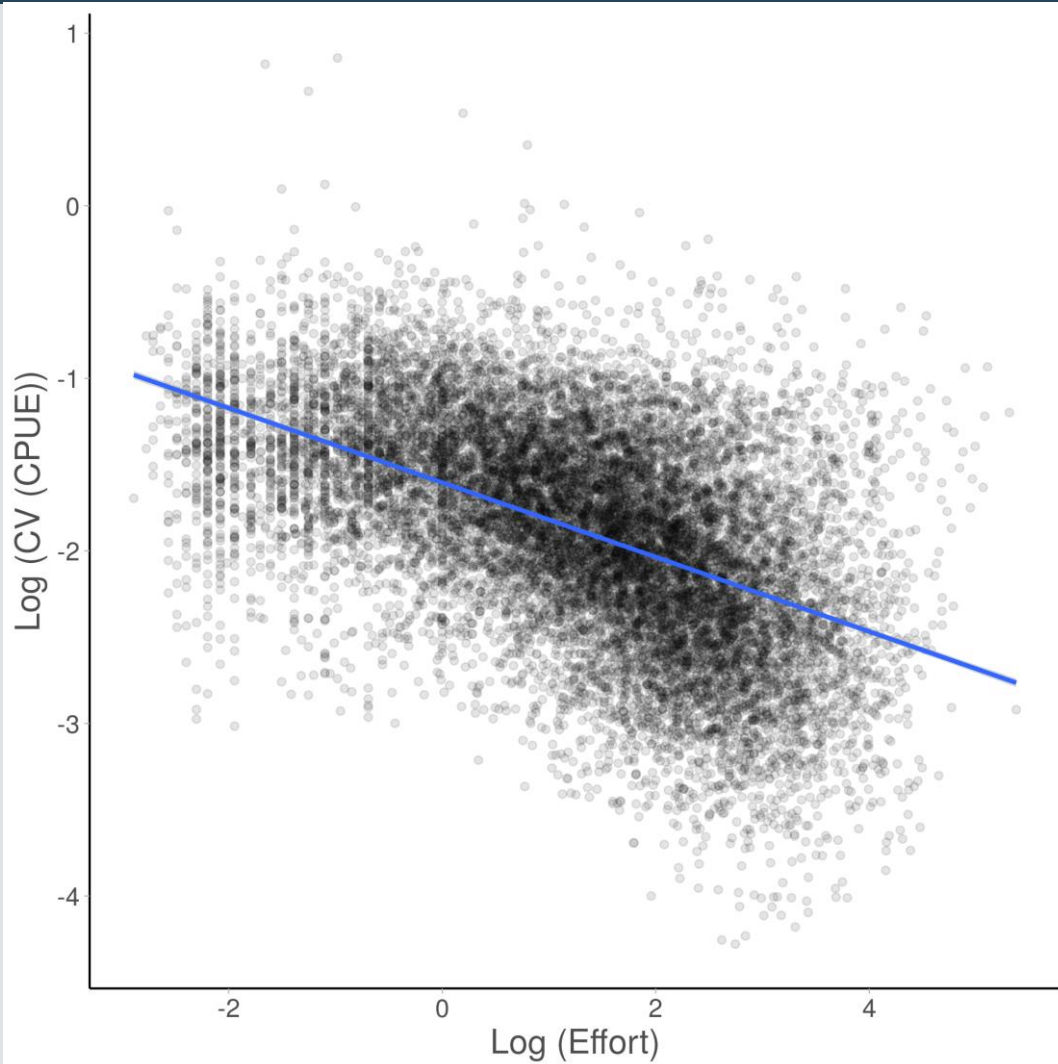


Predicted CV(CPUE)



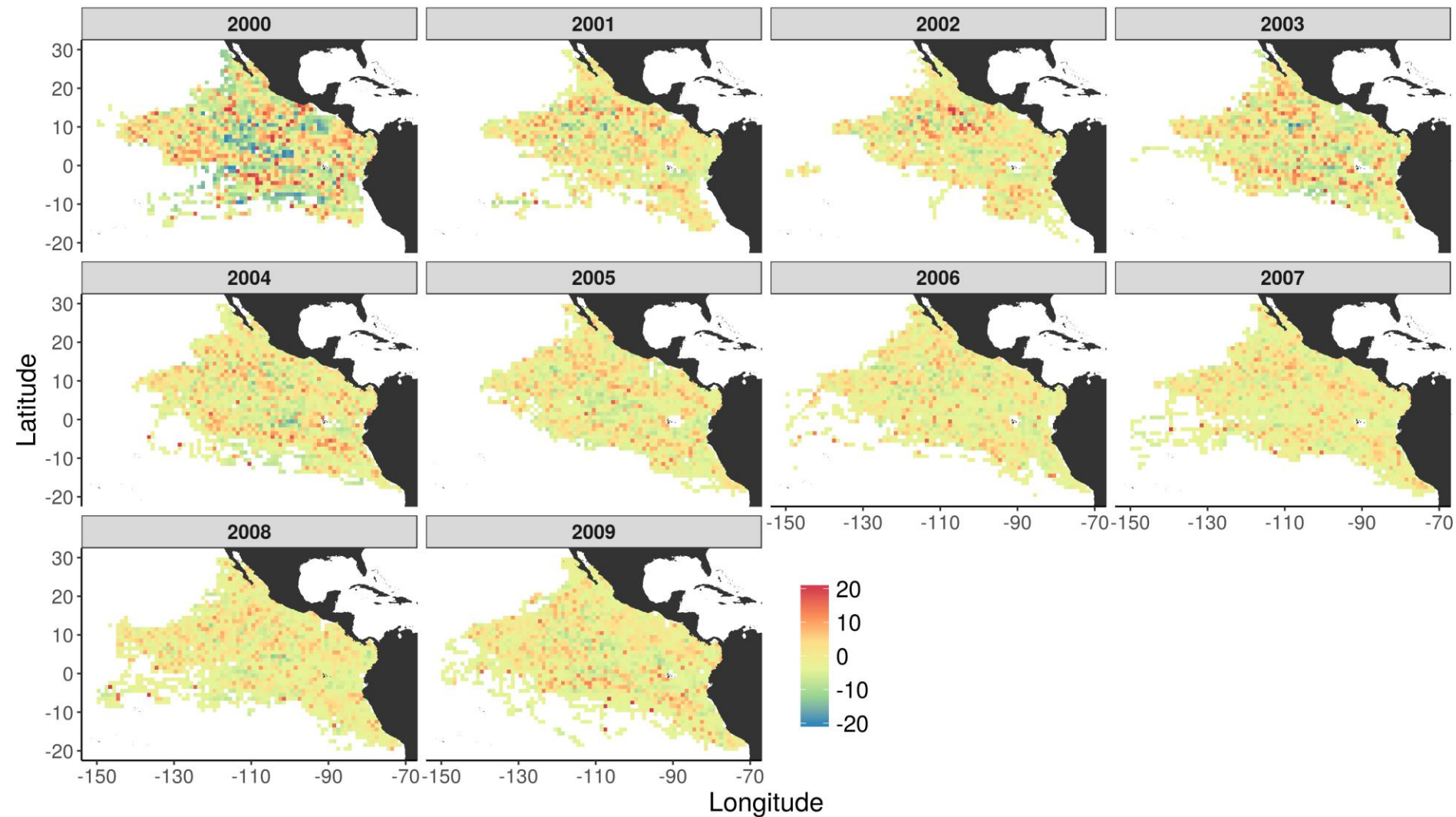
3. Results

$\log(\text{CV})$ vs. $\log(\text{Effort})$



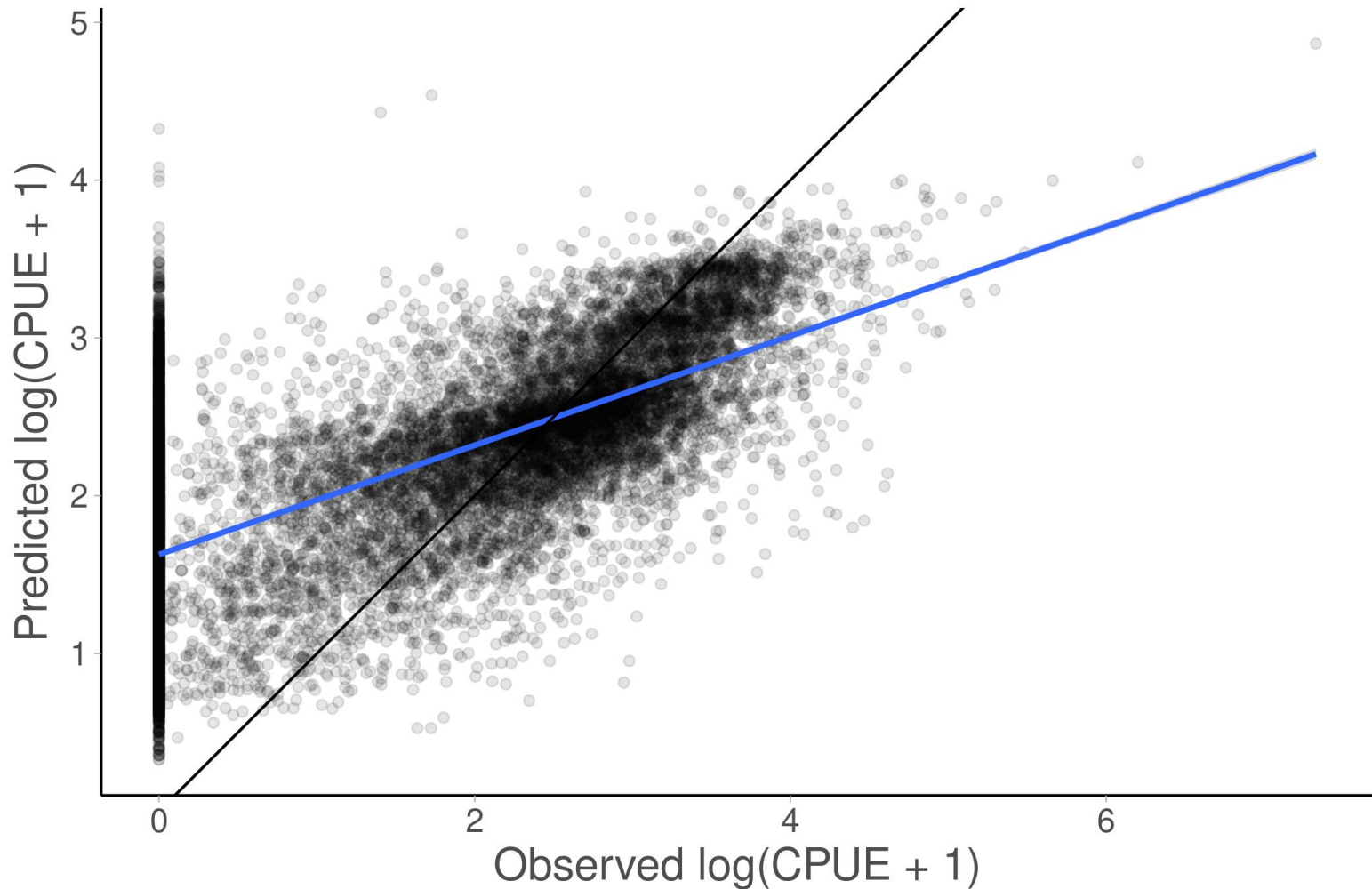
3. Diagnostics

Standardized residuals



3. Diagnostics

Bias (regression to the mean)




3. Diagnostics

Uncertainty estimates

Need *covariance* between spatiotemporal predictions

Rapidly evolving... 34,336 citations Breiman (2001)

1. Quantile regression forests – prediction quantiles
(‘ranger’, ‘grf’, Meinshausen 2006)
2. Jackknife & infinitesimal jackknife – standard error
(‘ranger’, Wager et al. 2014) 
3. U-statistics – asymptotically normal variance estimate
(‘surfin’, Mentch & Hooker 2016)
4. Generalized random forests – asymp. normal variance est. 
(‘grf’, Athey et al. 2017)
5. Bayesian additive regression trees – full posterior inference
(‘bayesMachine’, ‘dbarts’, ‘BART’, Chipman et al. 2010)

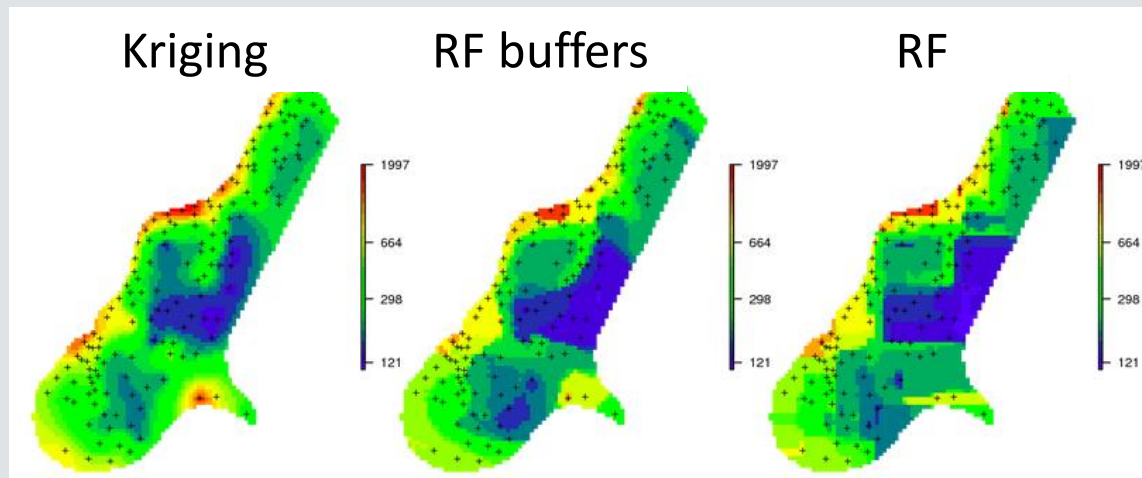
Other thoughts

Multivariate response:

- Include model.matrix as covariates:

```
levels(Data_Geostat$spp) <- c("A_stomias", "G_chalcogrammus", "H_lassodon")
sp.mat <- data.frame(model.matrix(~ spp - 1, Data_Geostat))
mv.dat <- cbind(Data_Geostat, sp.mat)
rfmv = ranger(Catch_KG ~ Lat + Lon + Year + sppA_stomias + sppG_chalcogrammus + sppH_lassodon,
              data=mv.dat, num.trees=1000, mtry=2, keep.inbag=T, write.forest=T)
```

Buffer distances to smooth predictions:



What we want (from Rick Methot)

- ★ Fast (coding vs. runtime vs. interpretation)
- ★ Replicable (method well-defined, get same answer)
- ★ Robust (insensitive to distributional assumptions, outliers)
- ★ Predictive ability (minimal errors, fill in space/time gaps)
- ★ Covariate effects (nonlinear, interactions)
- Uncertainty estimates (with known properties)
- ⊘ Specifiable structure (e.g. correlation through time, biology)
- ⊘ Unbiased (relative vs. absolute abundance)

Thank you!

SIO

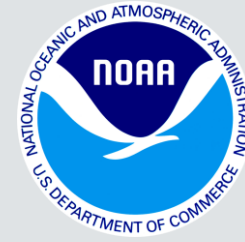
- Brice Semmens

SWFSC

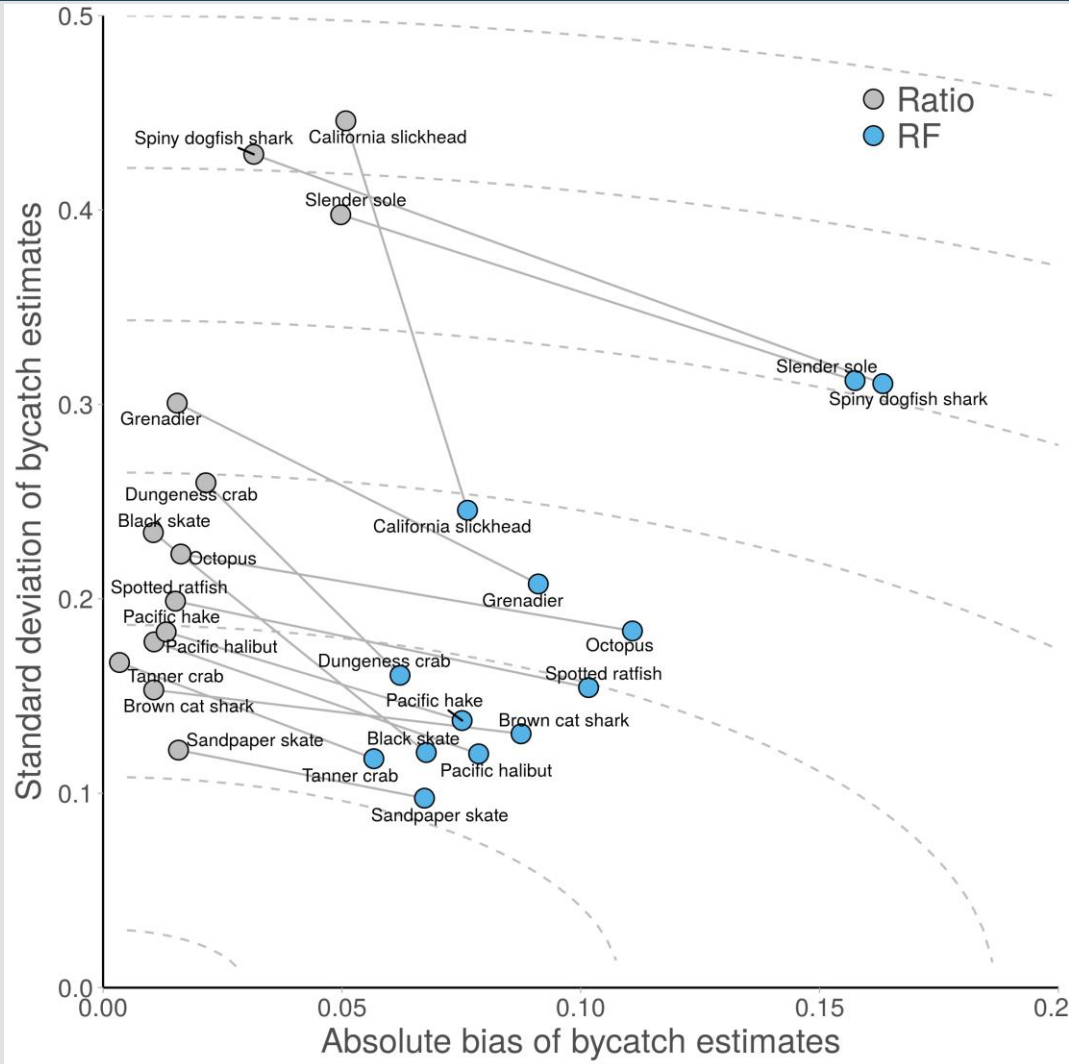
- Tomo Eguchi

NWFSC

- Eric Ward
- Jim Thorson
- Essential Fish Habitat (Blake Feist)
- West Coast Groundfish Observer Program (Jason Jannot)

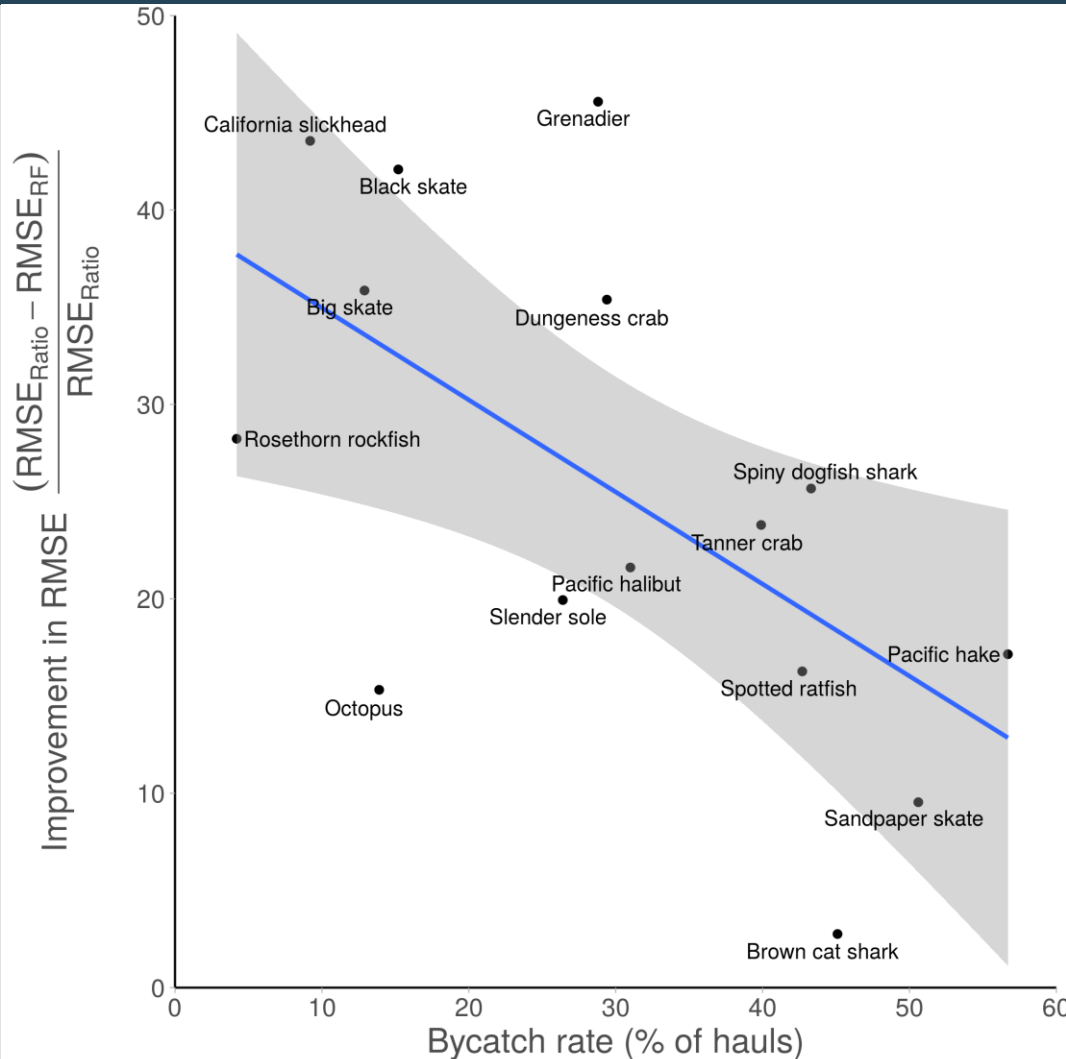


Bias-variance tradeoff by species...



2. Results

More worthwhile for rarer species



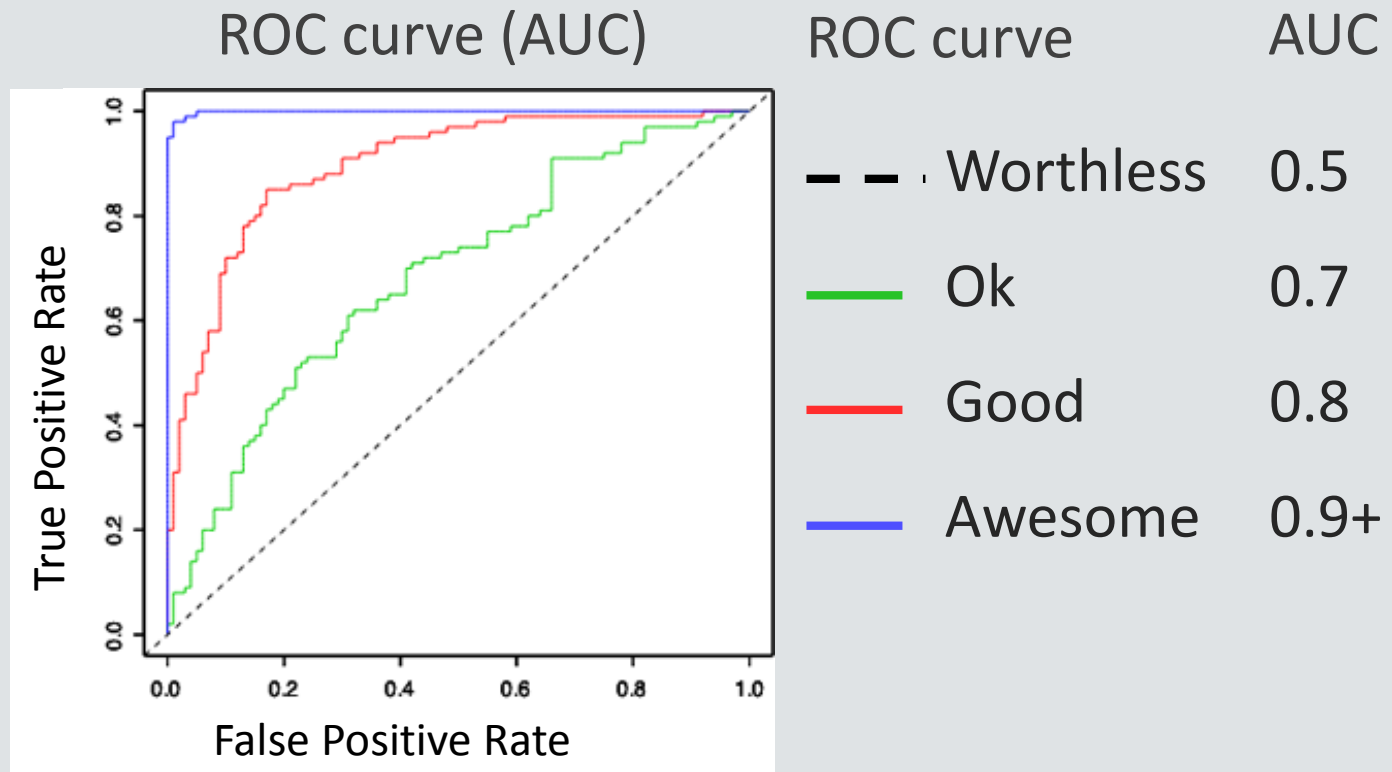
2. Results

Q1: Which spatial model is best?

Goal: *prediction*

5-fold cross validation repeated 10x

Binomial



Methods: evaluation

Q1: Which spatial model is best?

Goal: *prediction*

5-fold cross validation repeated 10x

Binomial

AUC

Positive

RMSE, R^2 (pred – obs)

$$\sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{n}}$$

West Coast Groundfish covariates

Binomial

Positive

~ sst + sst² +
depth + depth² +
distance to rocky substrate +
size of rocky patch +
in Rockfish Conservation Area +
predicted occurrence (survey) +
day of year +
spatial field

Hawaii Longline covariates

Binomial

Positive

~ sst + sst² +
day of year +
spatial field

RF

- + Better at prediction
- + More complex covariate relationships (incl. interactions)
- + Easier to set up and run
- + Not just a “black box”?

GMRF

- + Statistical inference, marginal posteriors for covariate effects
- + Ability to include observation error

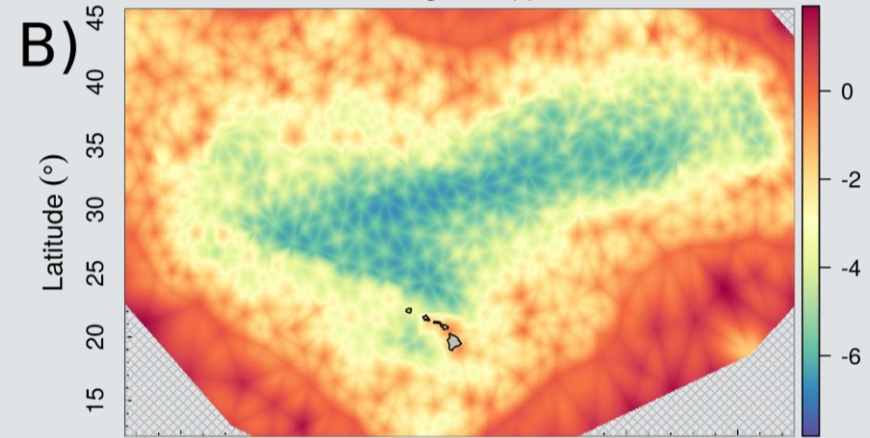
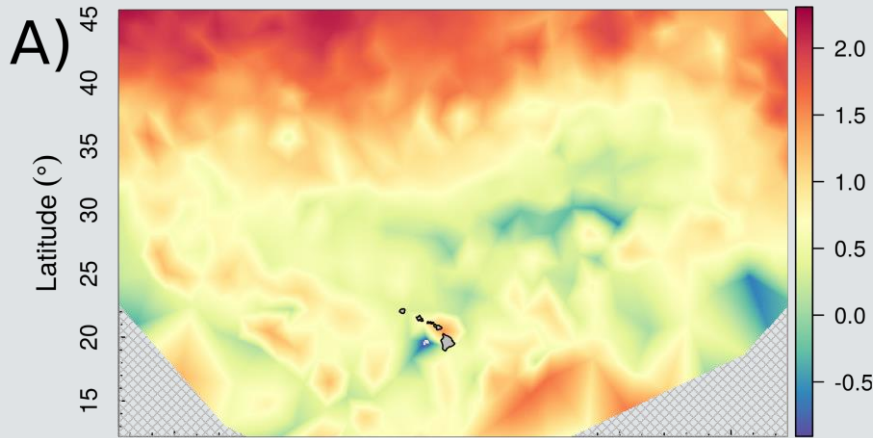
Variance of predictions



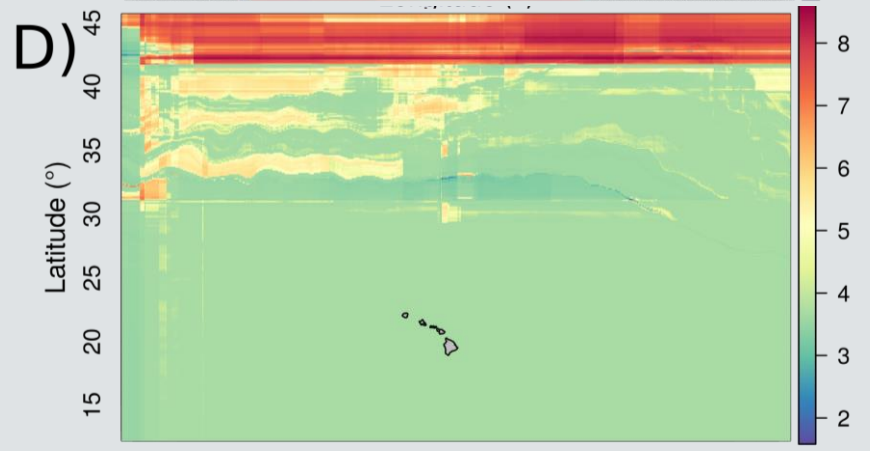
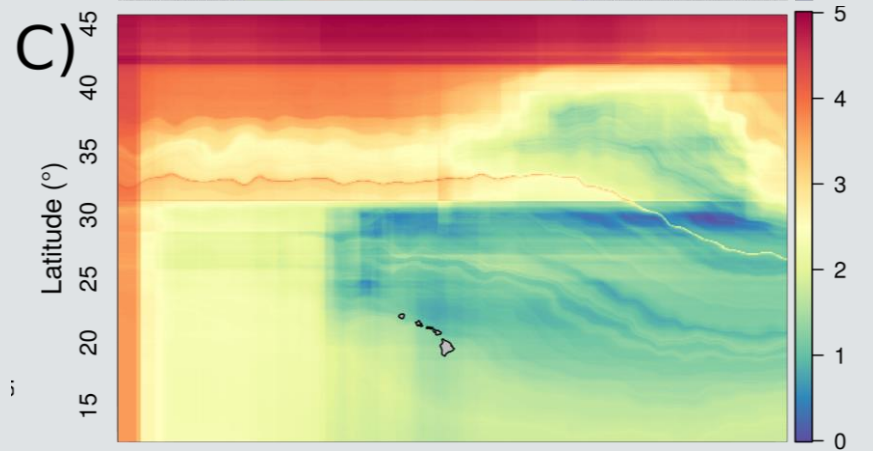
Mean

Variance

GMRF



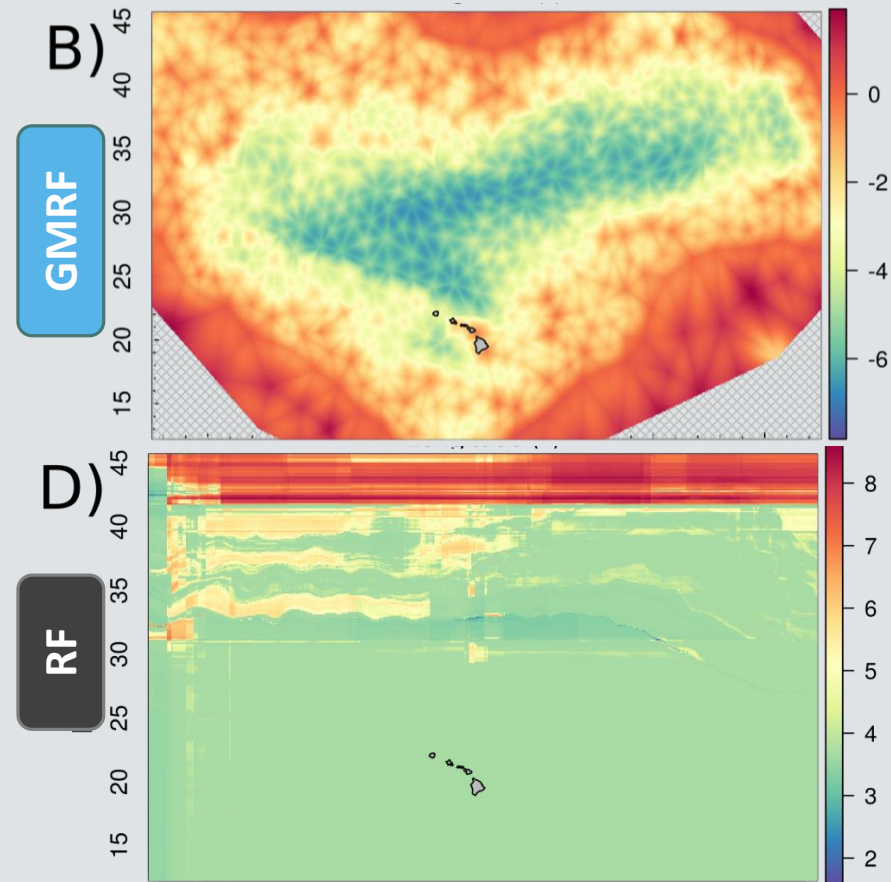
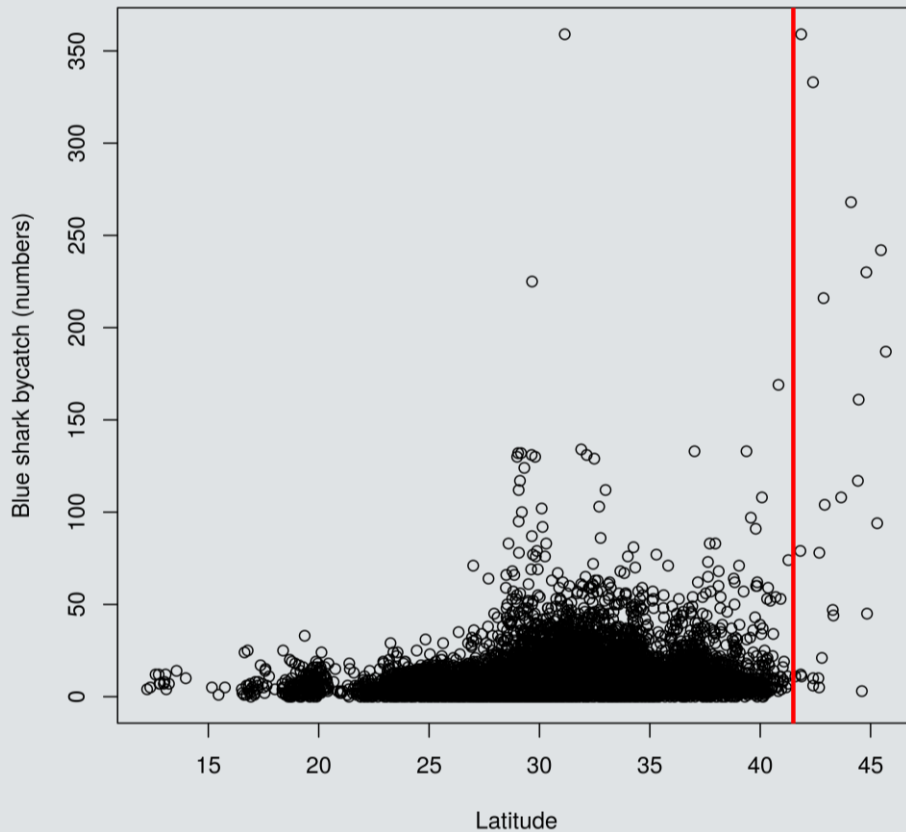
RF



Variance of predictions



Variance

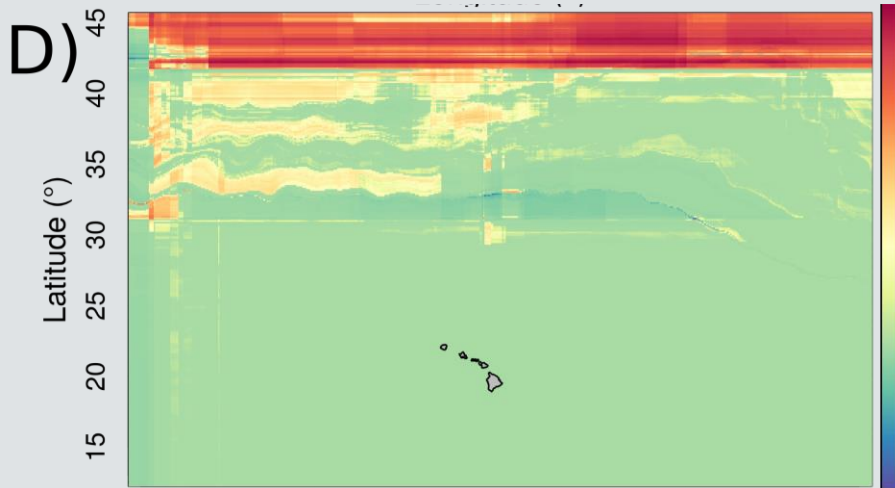


Variance of predictions

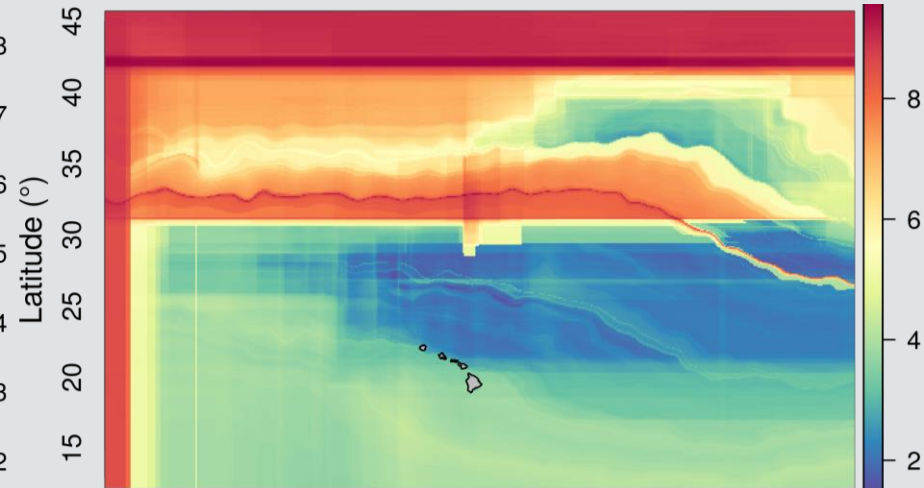
RF



Variance



Var(ind trees)



Non-parametric delta method /
“infinitesimal jackknife”