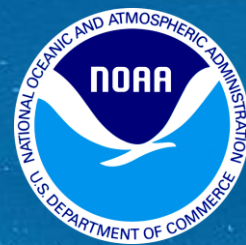


Photo: Norbert Wu Productions



What spatial statistical model is best for predicting fisheries bycatch risk?

BRIAN STOCK

Thank you!

SIO

- Brice Semmens

SWFSC

- Tomo Eguchi

NWFSC

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

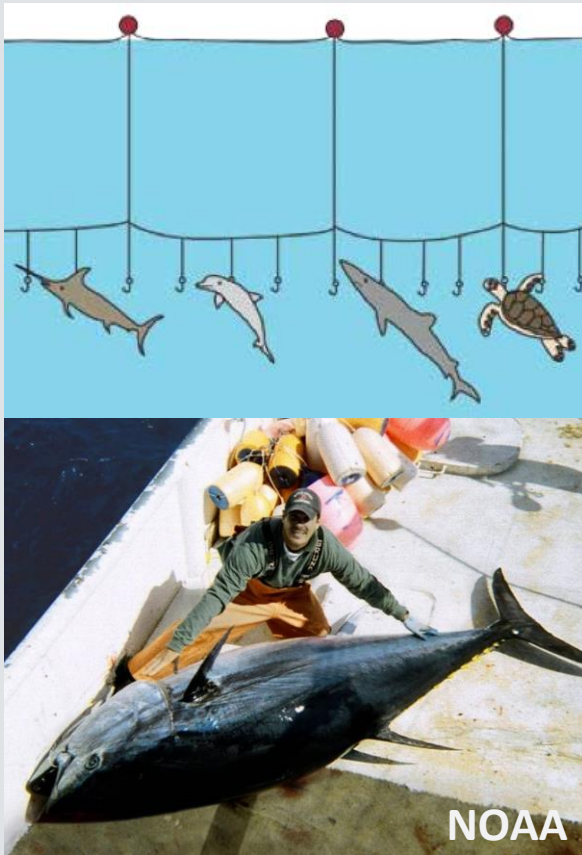
PIFSC

- Hawaii Longline Observer Program (Eric Forney)



“Target” vs. “bycatch”

Longline

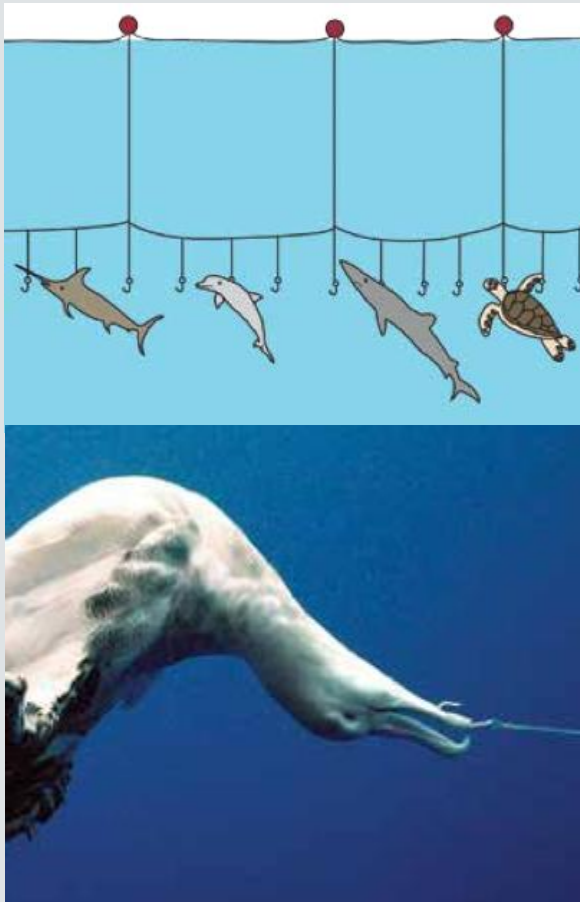


Trawl



“Target” vs. “bycatch”

Longline



Trawl



Bycatch is a big (spatial) issue

Protected species



Recovering species



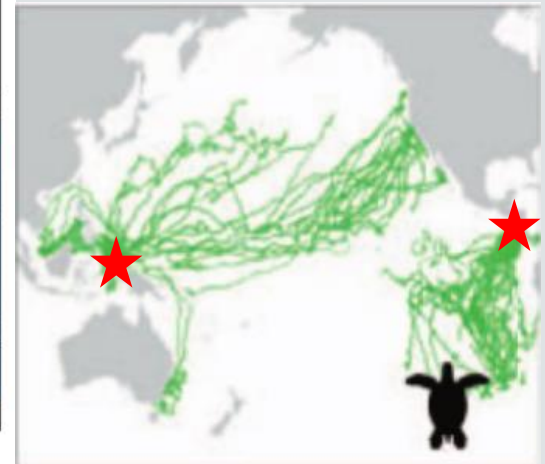
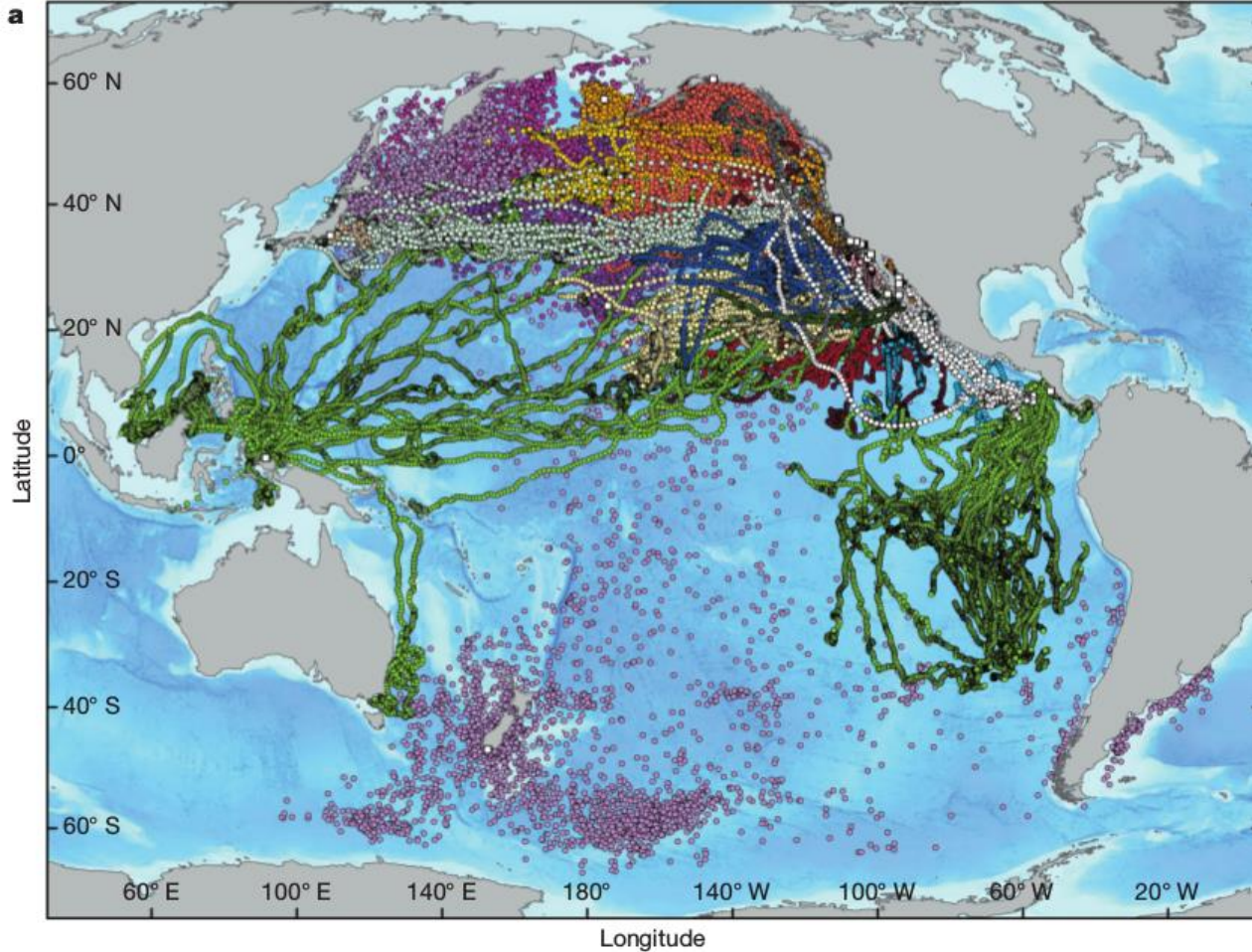
Competing fisheries



Unmarketable species



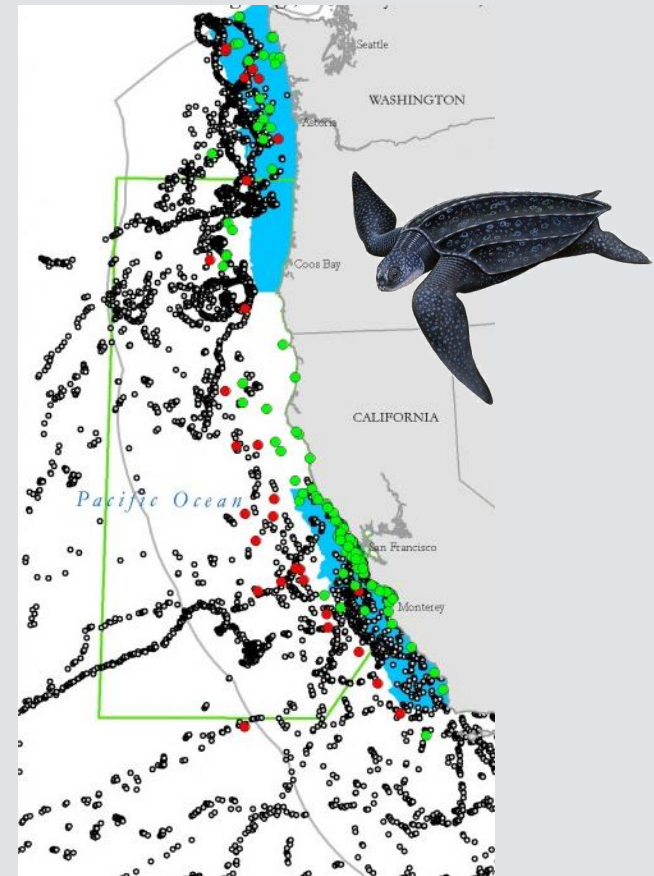
Difficult when they move so much...



Static vs. dynamic management

Dynamic

Static

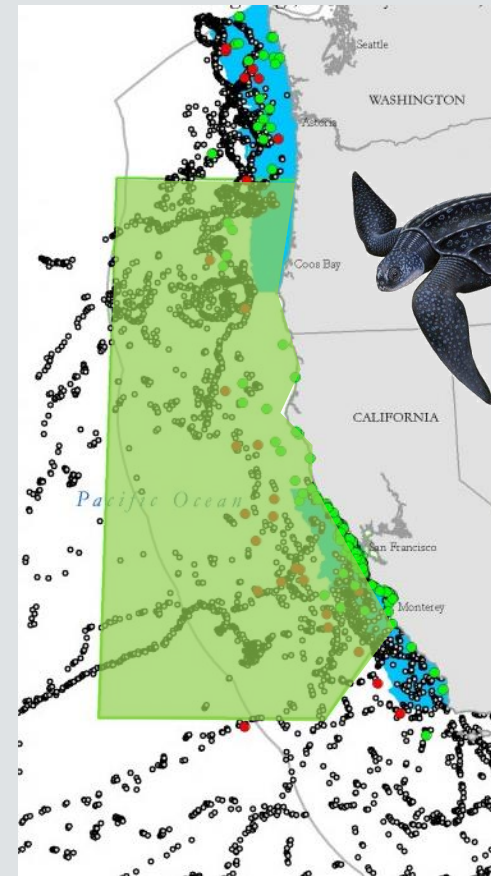


Static vs. dynamic management

Dynamic

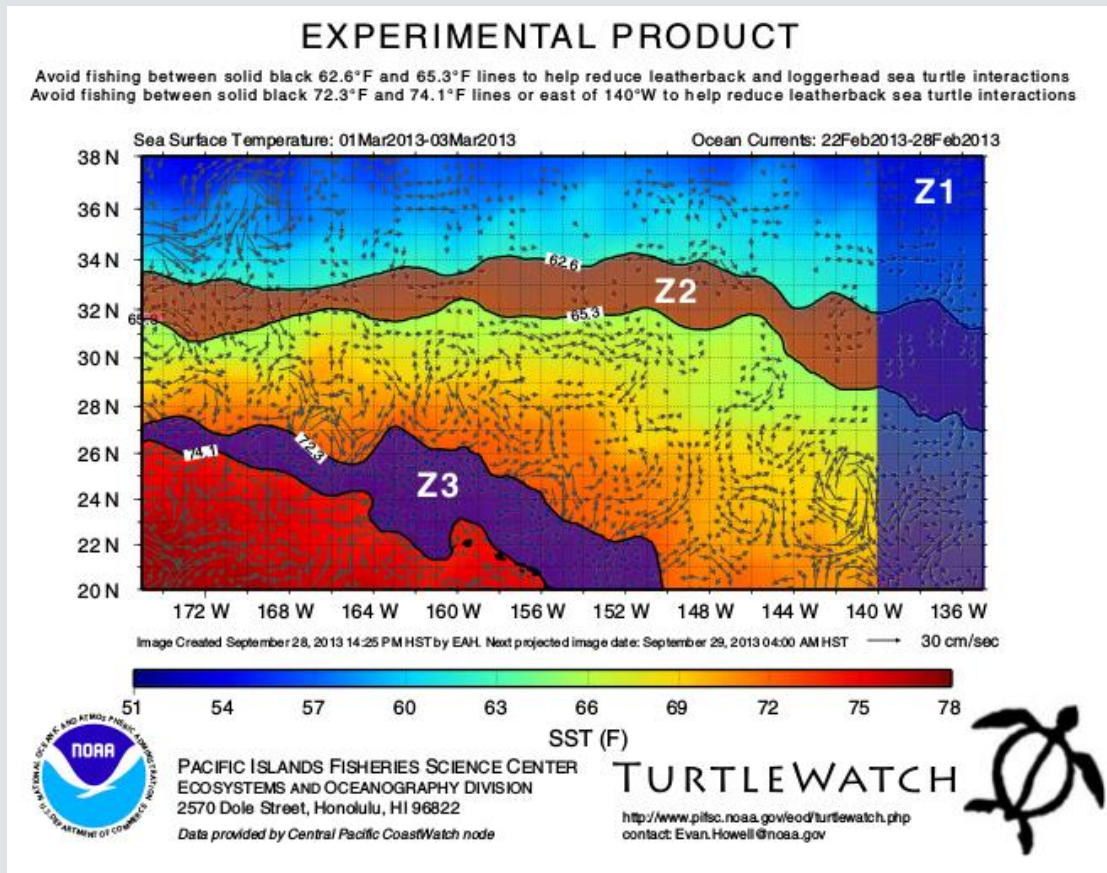
1. Effectively protected?
2. Huge loss of fishing area

Static

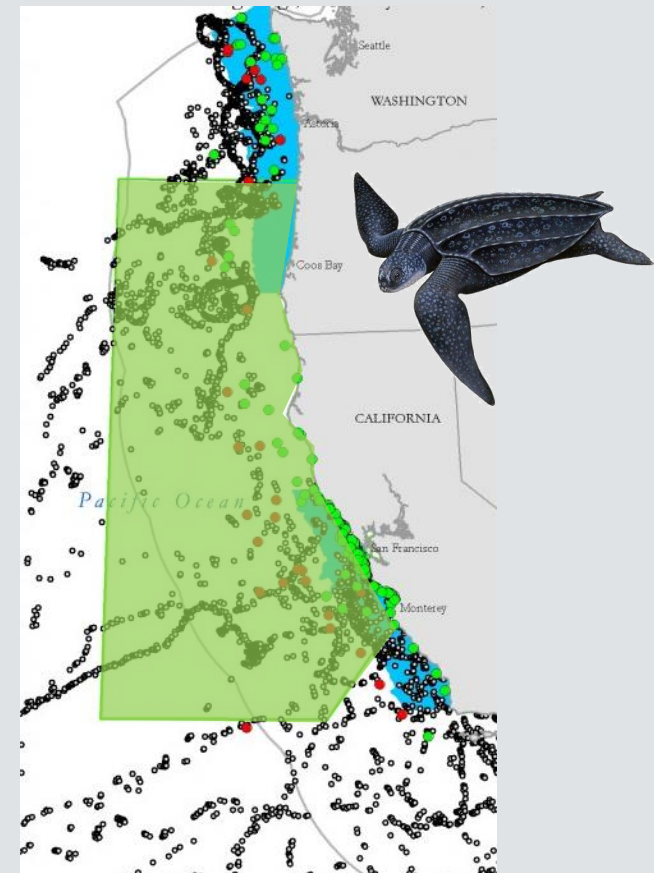


Static vs. dynamic management

Dynamic

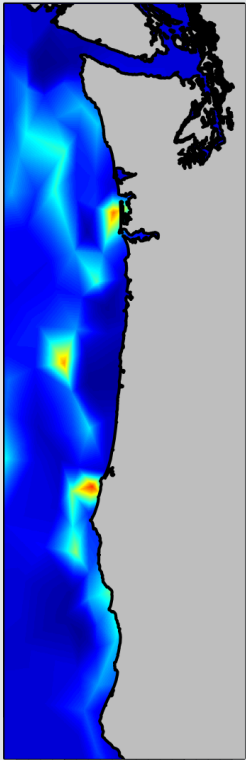


Static



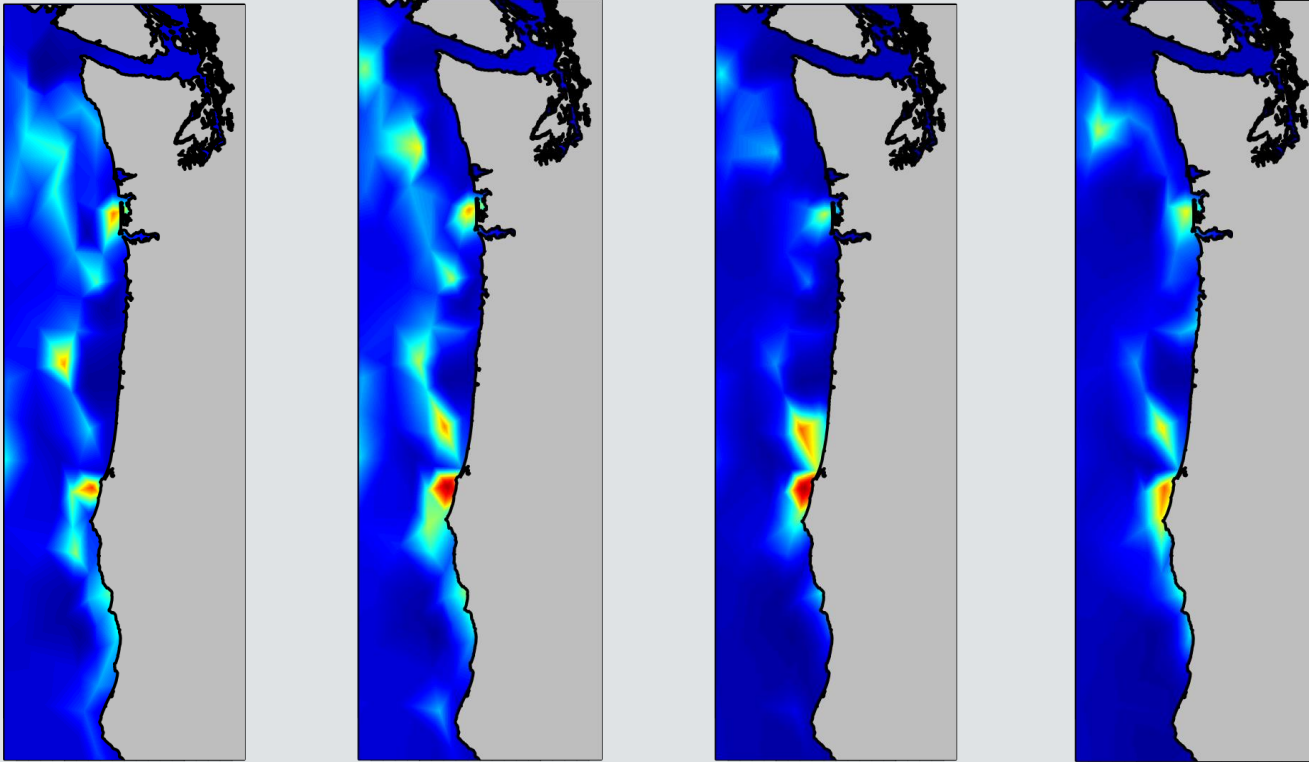
Tools for dynamic management

Need map of bycatch risk



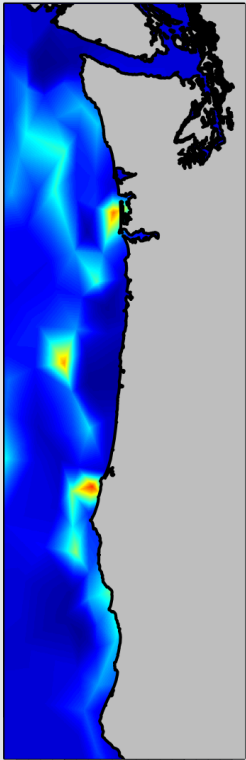
Tools for dynamic management

Need map of bycatch risk



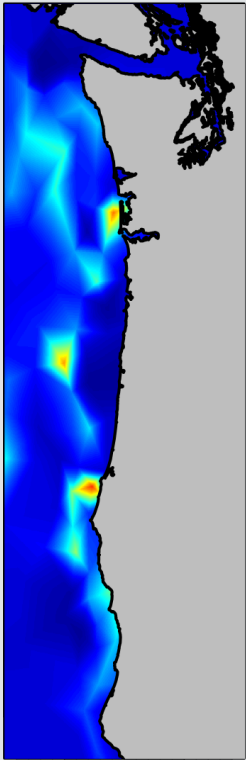
Tools for dynamic management

Need map of bycatch risk



- temperature
- depth
- substrate
- spatial field

Q1: Which spatial model is best?



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

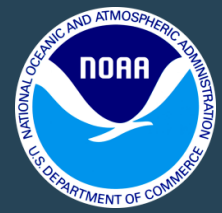
GLM

GAM

GMRF

RF

The data (fisheries observers)



West Coast Groundfish Trawl

- 2002-2013
- 55,835 tows



Hawaii Longline

- 1994-2014
- 16,714 sets (swordfish only)



Research question

Q2: Does the answer depend on species?



Habitat: Bottom

Movement: Med

Bycatch Rate: 29%



Habitat: Rocky bottom

Movement: Low

Bycatch Rate: 18%



Habitat: Rocky bottom

Movement: Low

Bycatch Rate: 0.3%



Habitat: Open ocean

Movement: High

Bycatch Rate: 96%



Habitat: Open ocean

Movement: High

Bycatch Rate: 1.4%



Habitat: Open ocean

Movement: High

Bycatch Rate: 0.7%

Research question

Q2: Does the answer depend on species?



Habitat: Bottom

Movement: Med

Bycatch Rate: 29%



Habitat: Rocky bottom

Movement: Low

Bycatch Rate: 18%



Habitat: Rocky bottom

Movement: Low

Bycatch Rate: 0.3%



Habitat: Open ocean

Movement: High

Bycatch Rate: 96%



Habitat: Open ocean

Movement: High

Bycatch Rate: 1.4%



Habitat: Open ocean

Movement: High

Bycatch Rate: 0.7%

Q2: Does the answer depend on species?



Habitat:

Bottom

Rocky bottom

Rocky bottom

Movement:

Med

Low

Low

Bycatch Rate:

29%

18%

0.3%



Habitat:

Open ocean

Open ocean

Open ocean

Movement:

High

High

High

Bycatch Rate:

96%

1.4%

0.7%

Research question

Q2: Does the answer depend on species?



Habitat: Bottom

Movement: Med

Bycatch Rate: 29%



Habitat: Rocky bottom

Movement: Low

Bycatch Rate: 18%



Habitat: Rocky bottom

Movement: Low

Bycatch Rate: 0.3%



Habitat: Open ocean

Movement: High

Bycatch Rate: 96%



Habitat: Open ocean

Movement: High

Bycatch Rate: 1.4%

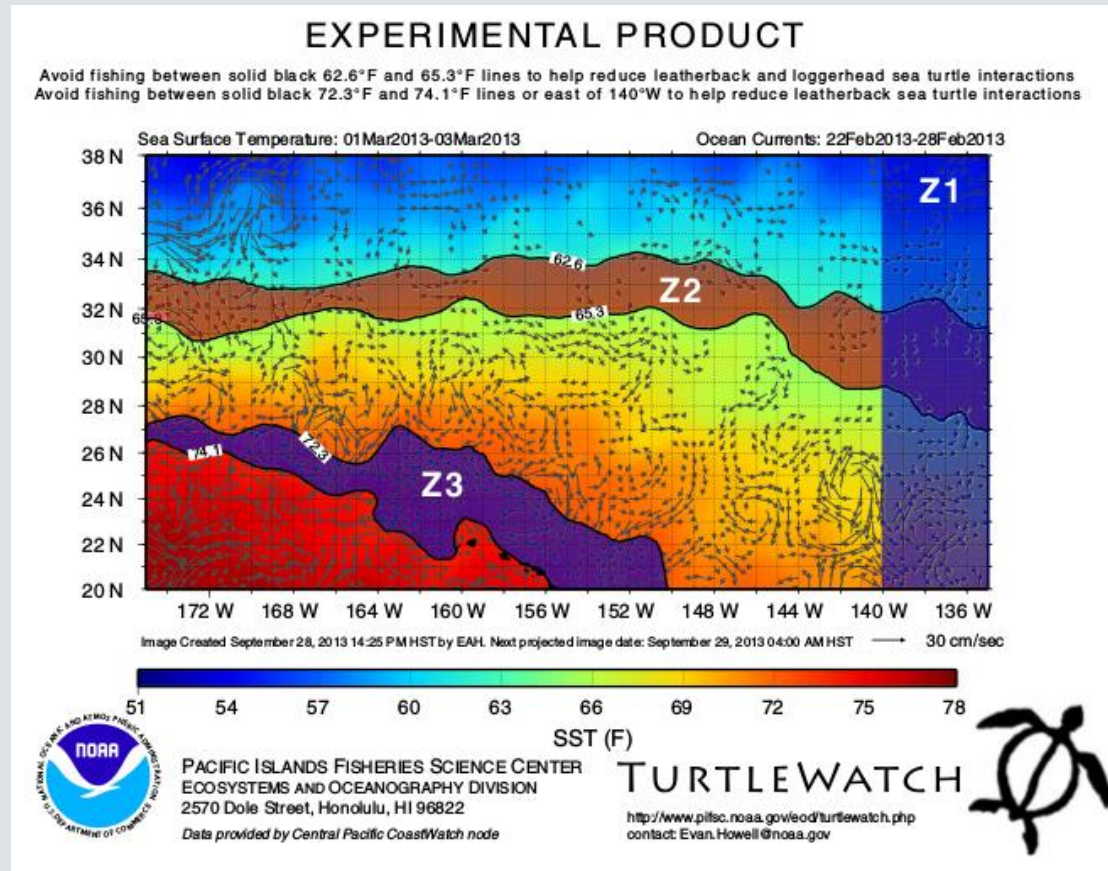


Habitat: Open ocean

Movement: High

Bycatch Rate: 0.7%

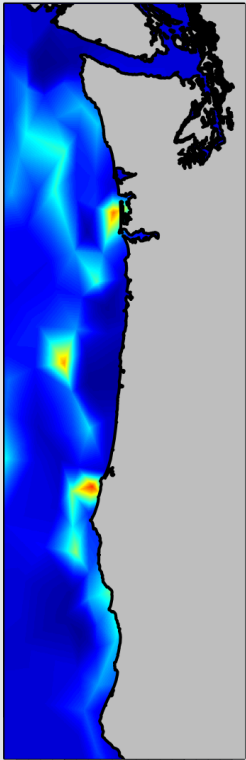
Q3: How much bycatch can they prevent?



Research question

“Species distribution models”

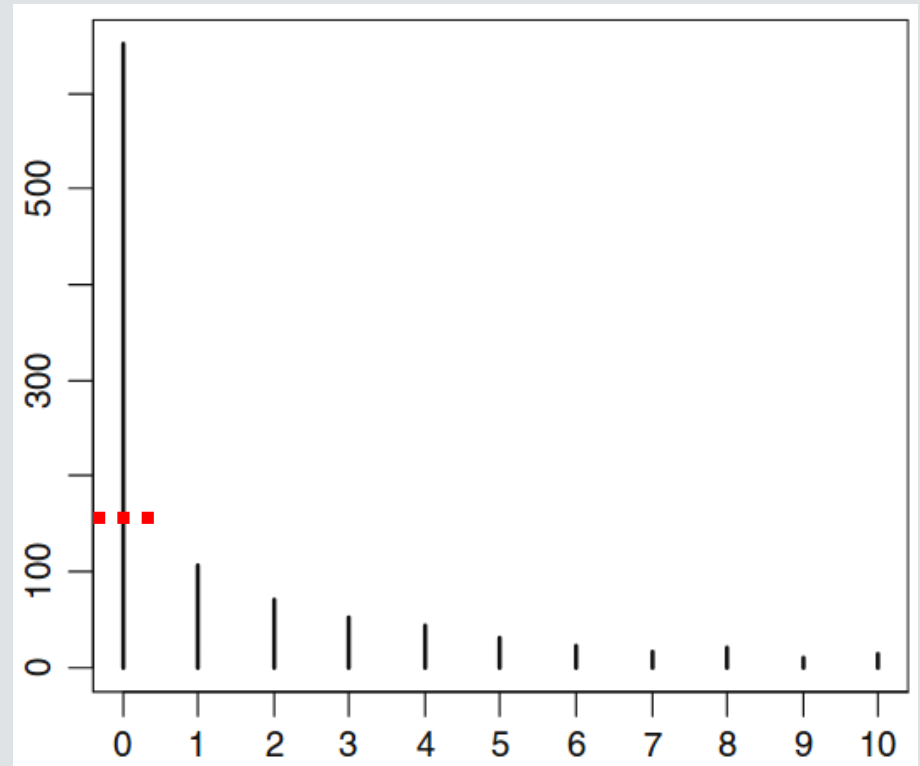
Fundamental ecological question: *where are they?*



- temperature
- depth
- substrate
- spatial field

“Zero-inflated” data

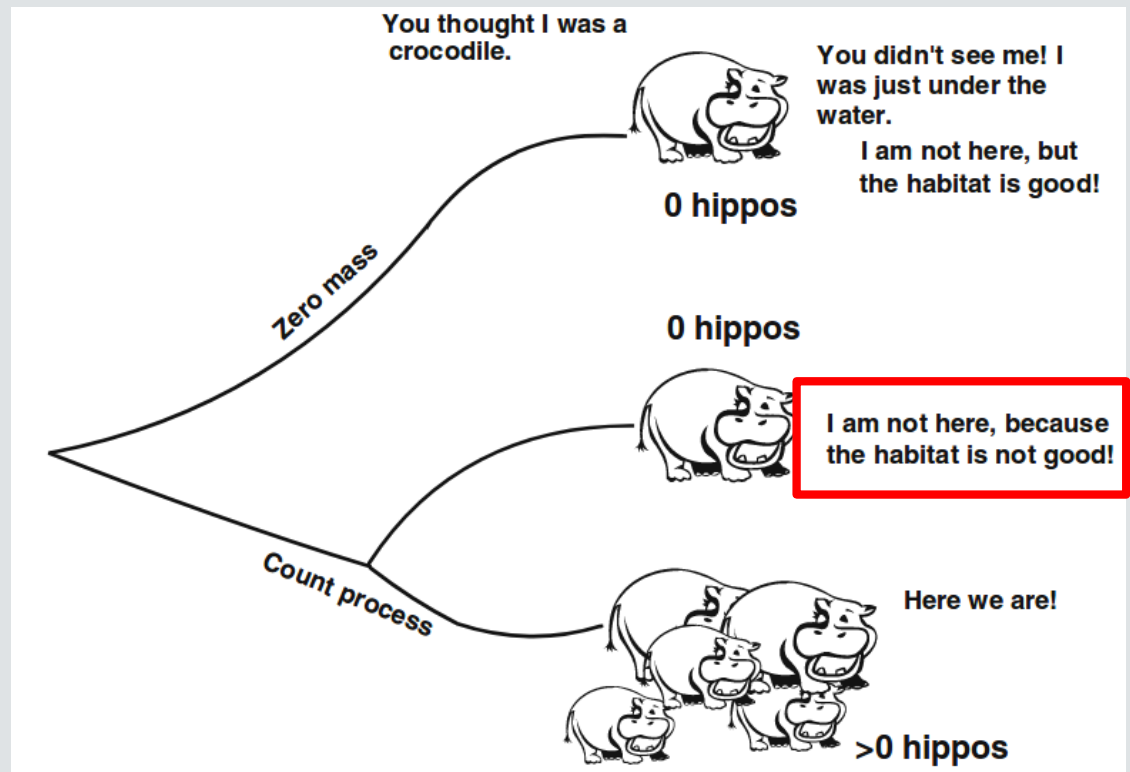
More zeros than expected



“Zero-inflated” data

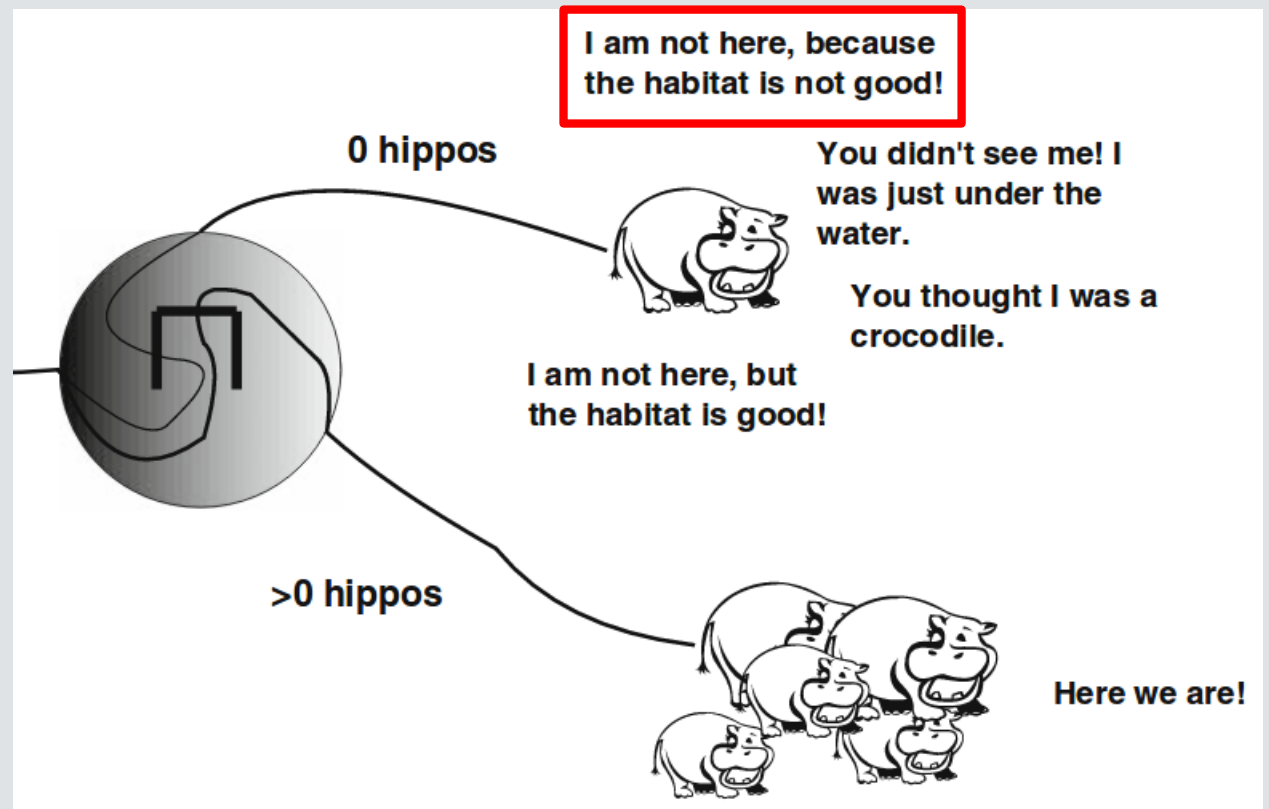
Approach 1: Zero-inflated distributions

- ZI-Poisson
- ZI-Neg Binomial



“Zero-inflated” data

Approach 2: Delta (hurdle) model



“Zero-inflated” data

Approach 2: Delta (hurdle) model

Binomial

Pr(some bycatch)

$$\text{logit}(p_i) = \log\left(\frac{p_i}{1 - p_i}\right) = \mathbf{X}_i\boldsymbol{\beta}$$

$$Y_i \sim \text{Bernoulli}(p_i)$$

“Zero-inflated” data

Approach 2: Delta (hurdle) model

Binomial

Pr(some bycatch)

Positive

E(bycatch | some bycatch)

$$\log(\mu_i) = \mathbf{X}_i\boldsymbol{\beta}$$

$$Y_i \sim \text{Gamma}(\mu_i, \nu) \quad \text{for } Y_i > 0$$

“Zero-inflated” data

Approach 2: Delta (hurdle) model

Binomial

$\Pr(\text{some bycatch})$

Positive

$E(\text{bycatch} \mid \text{some bycatch})$

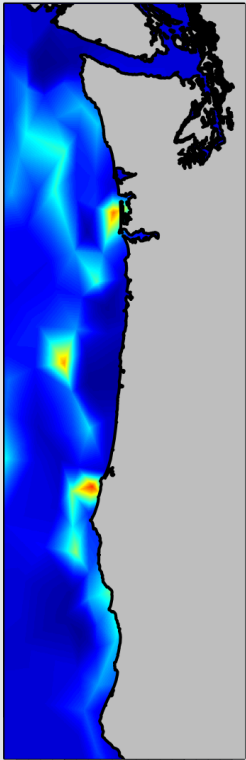
Binomial

x

Positive

$E(\text{bycatch})$

Q1: Which spatial model is best?



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

GLM

GAM

GMRF

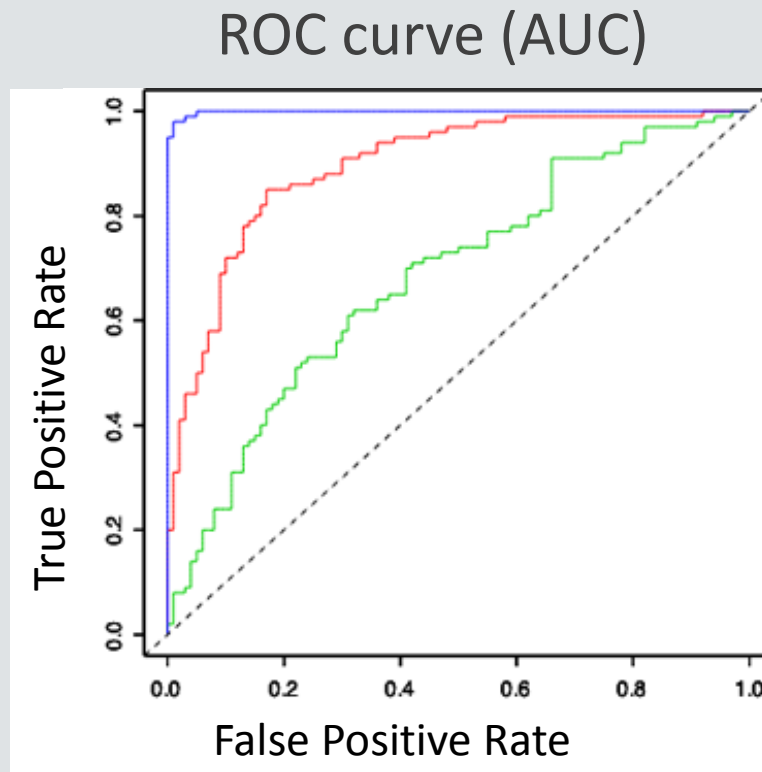
RF

Q1: Which spatial model is best?

Goal: prediction

5-fold cross validation repeated 10x

Binomial



ROC curve	AUC
--- Worthless	0.5
— Ok	0.7
— Good	0.8
— Awesome	0.9+

Methods: evaluation

Q1: Which spatial model is best?

Goal: prediction

5-fold cross validation repeated 10x

Binomial

AUC

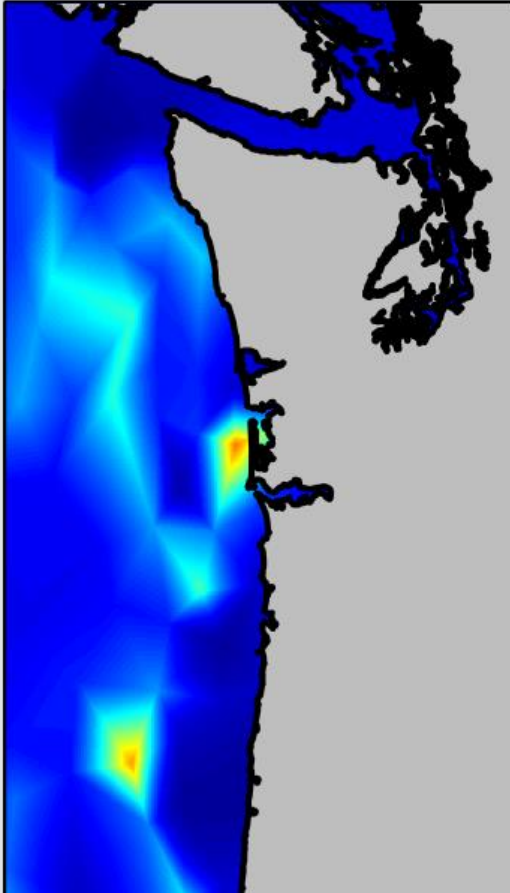
Positive

RMSE, R^2 (pred – obs)

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

Methods: evaluation

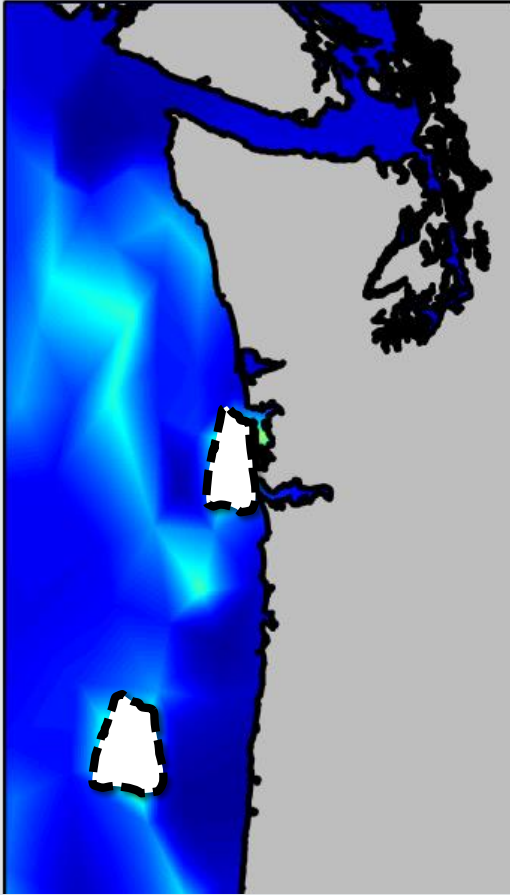
Q3: How much bycatch can they prevent?



Simulate management:

1. Predict bycatch risk at test locations

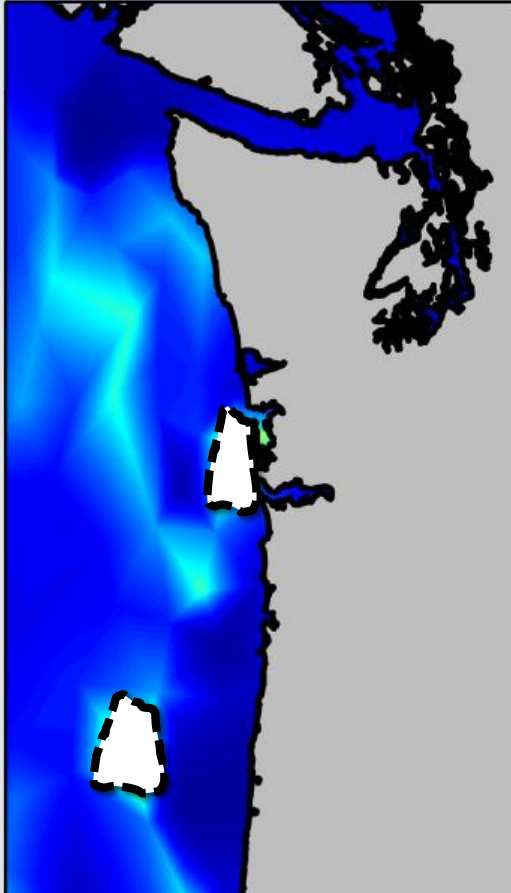
Q3: How much bycatch can they prevent?



Simulate management:

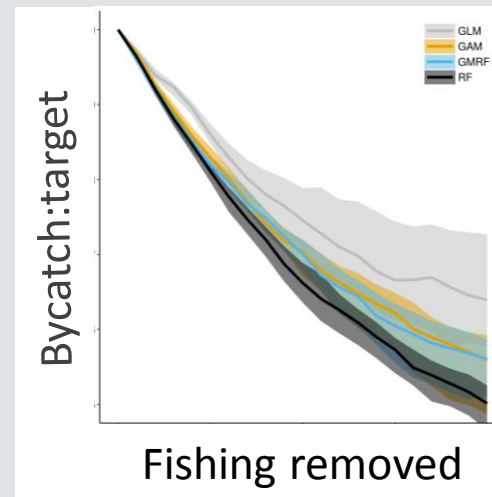
1. Predict bycatch risk at test locations
2. Remove X% of fishing effort with highest bycatch risk

Q3: How much bycatch can they prevent?



Simulate management:

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)



Q1: Which spatial model is best?

GLM

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

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

Binomial

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

Positive

GAM

GMRF

RF

Q1: Which spatial model is best?

GLM

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

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

Binomial

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

Positive

GAM

GMRF

RF

How much variability can we explain?

- with covariates
- *without spatial locations*

Q1: Which spatial model is best?

GLM

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

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

Binomial

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

Positive

GAM

GMRF

RF

Problem:

spatial correlation in residuals (Prediction – Observed)

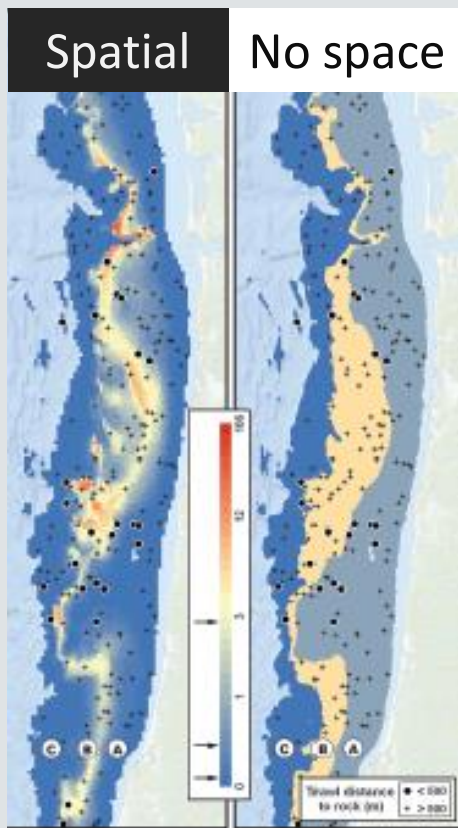


Why does spatial correlation matter?

1. Valid statistical inference
 - Observations not independent
 - Lower effective sample size (i.e. CI should be wider)

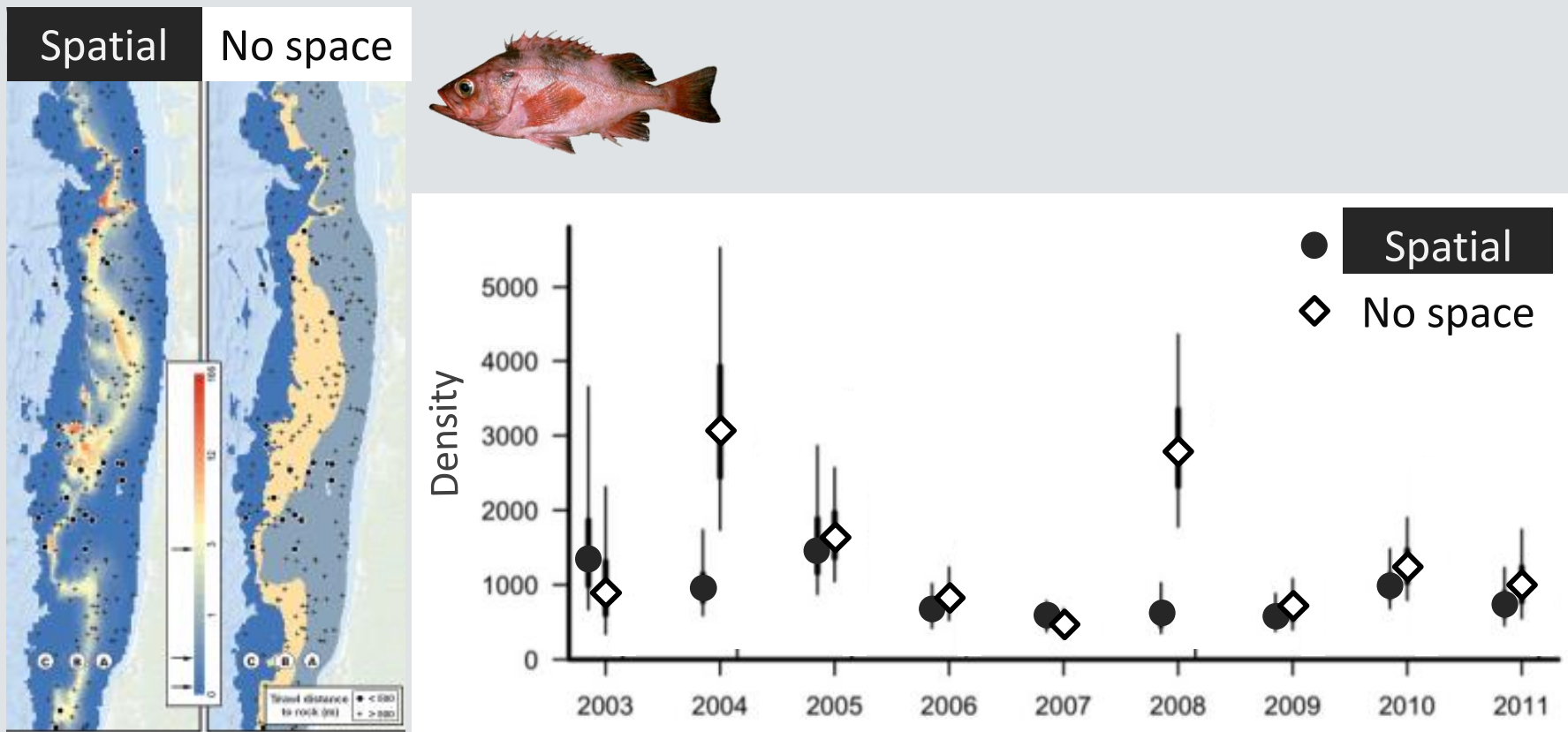
Why does spatial correlation matter?

2. Get the temporal trend right



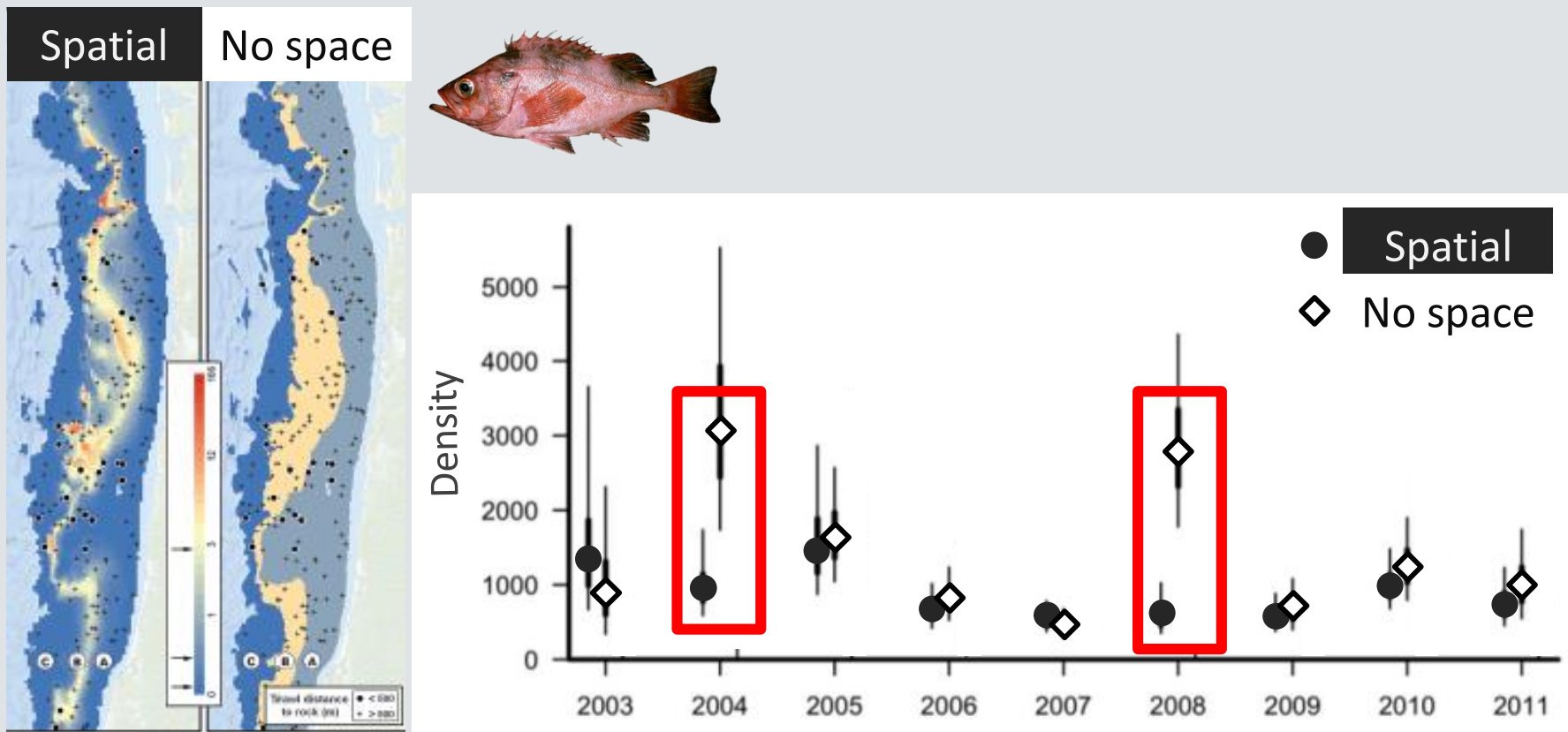
Why does spatial correlation matter?

2. Get the temporal trend right



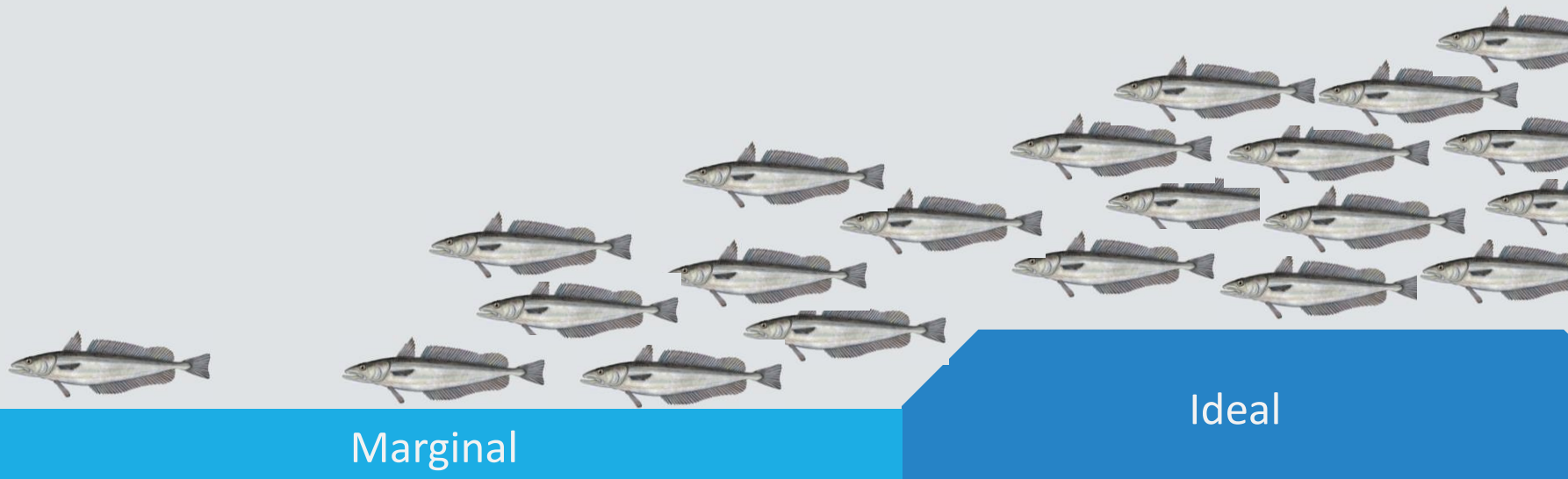
Why does spatial correlation matter?

2. Get the temporal trend right



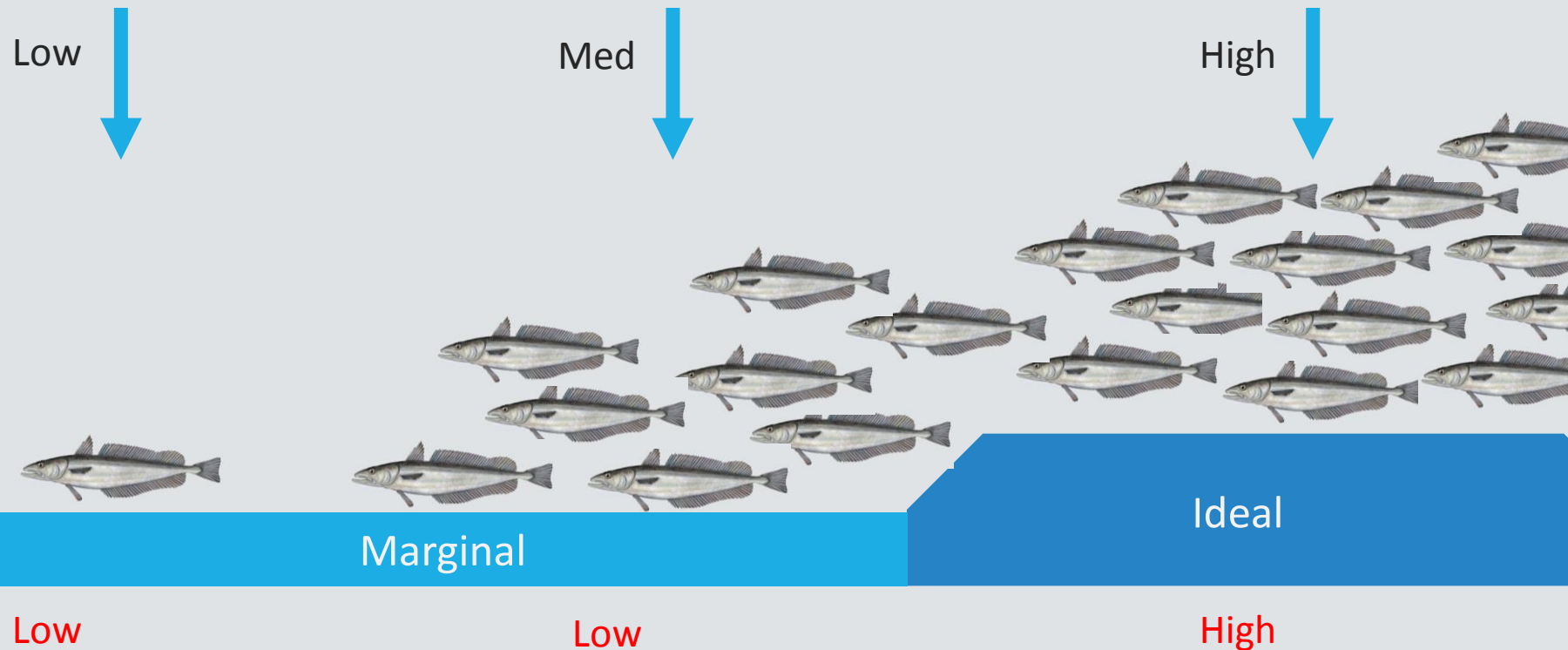
Why does spatial correlation matter?

3. Effect of habitat vs. schooling



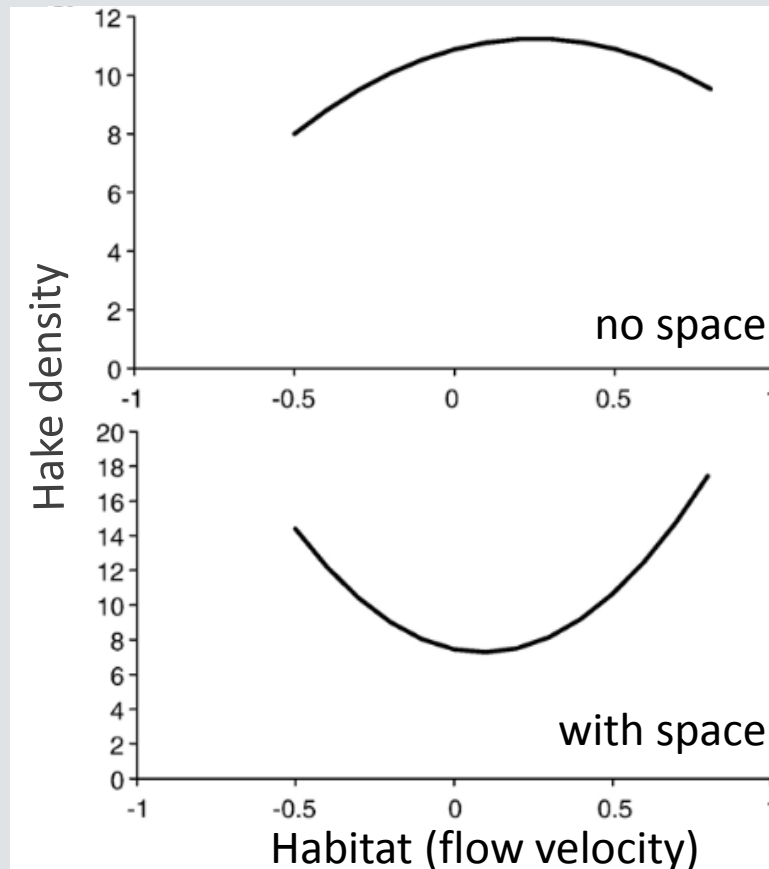
Why does spatial correlation matter?

3. Effect of habitat vs. schooling



Why does spatial correlation matter?

3. Effect of habitat vs. schooling



Q1: Which spatial model is best?

GLM

GAM

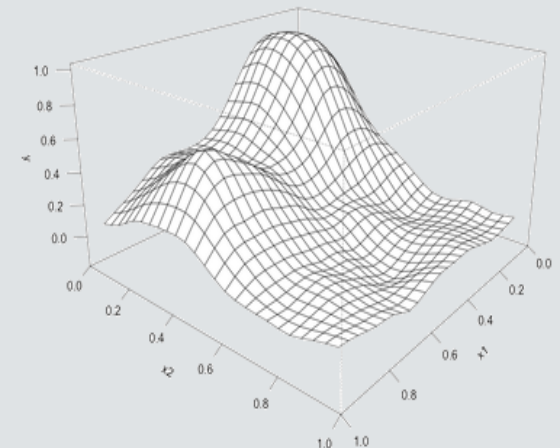
GMRF

RF

Generalized Additive Models

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

- Common, simple approach
- Parameterized by spline basis functions
(*not spatial correlation*)



Q1: Which spatial model is best?

GLM

Gaussian Markov random field

GAM

GMRF

RF

Q1: Which spatial model is best?

GLM

GAM

GMRF

RF

Gaussian Markov random field

- Models *covariance* as function of spatial locations
obs \sim environmental predictors + $MVN(0, \Sigma)$

Q1: Which spatial model is best?

GLM

GAM

GMRF

RF

Gaussian Markov random field

- Models *covariance* as function of spatial locations
obs \sim environmental predictors + $MVN(0, \Sigma)$

Problem...

- Σ has $O(N^2)$ elements
- Computations scale as $O(N^3)$ from $|\Sigma|$ and Σ^{-1}

Q1: Which spatial model is best?

GLM

GAM

GMRF

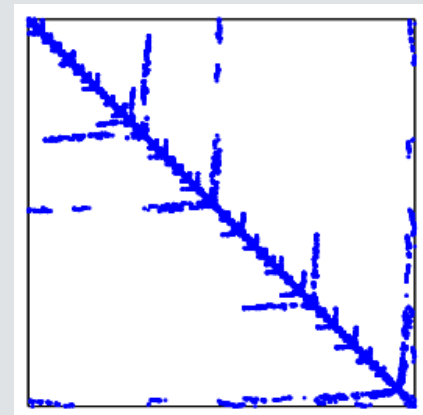
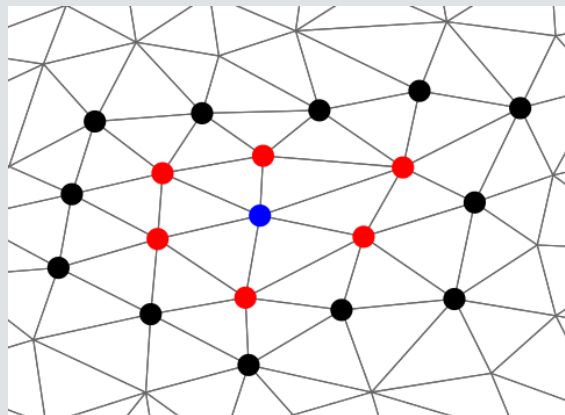
RF

Gaussian **Markov** random field

- Models *covariance* as function of spatial locations
obs \sim environmental predictors + $MVN(0, \Sigma)$

Solution:

- correlation = 0 for “far away” points \rightarrow sparse matrix



Q1: Which spatial model is best?

GLM

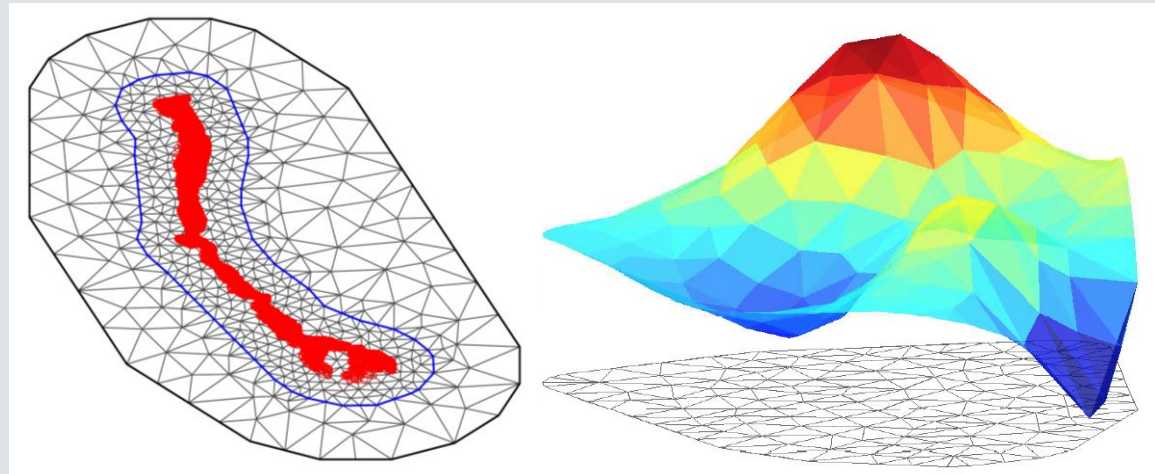
GAM

GMRF

RF

Gaussian Markov random field

- Models *covariance* as function of spatial locations
obs \sim environmental predictors + $MVN(0, \Sigma)$
- Discrete approximation of continuous space



Q1: Which spatial model is best?

GLM

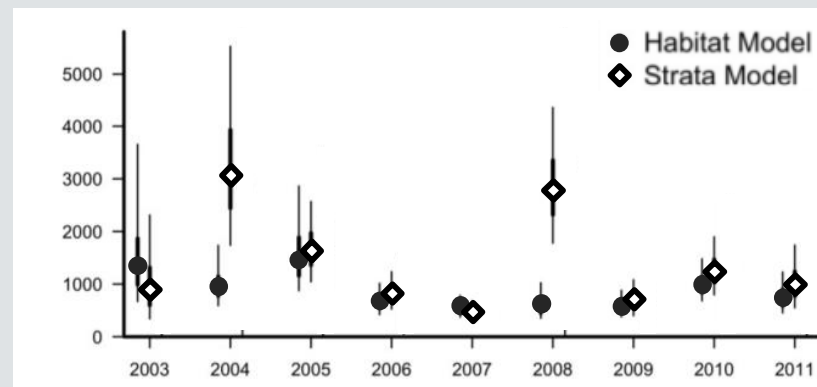
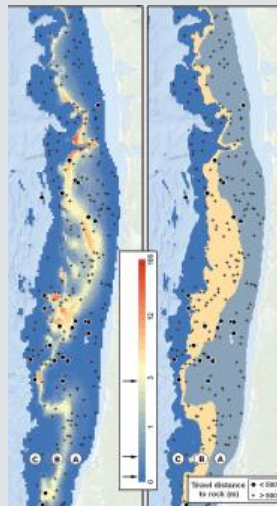
GAM

GMRF

RF

Gaussian Markov random field

- Models *covariance* as function of spatial locations
obs \sim environmental predictors + $MVN(0, \Sigma)$
- Increasing adoption in fisheries



Q1: Which spatial model is best?

GLM

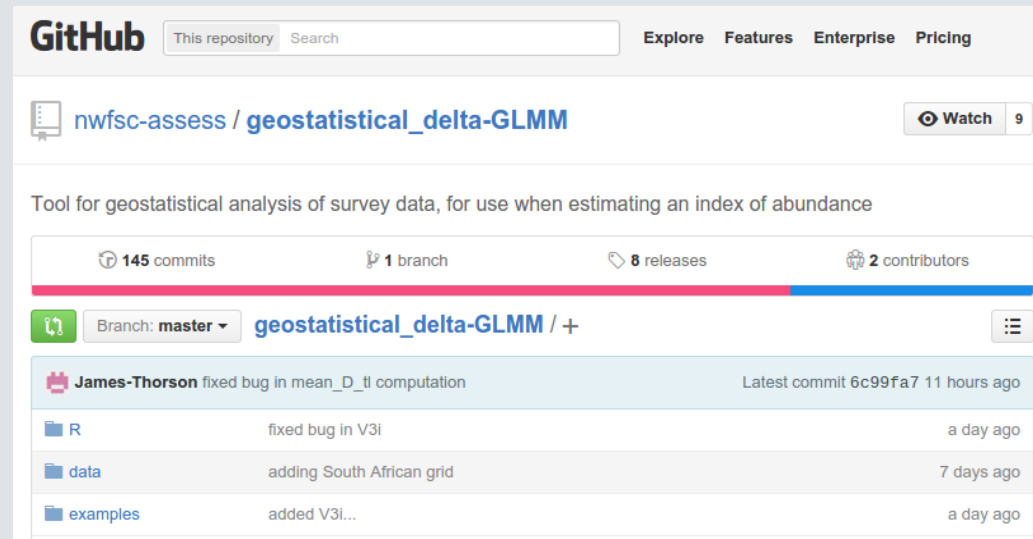
GAM

GMRF

RF

Gaussian Markov random field

- Models *covariance* as function of spatial locations
 $\text{obs} \sim \text{environmental predictors} + \text{MVN}(0, \Sigma)$
- Increasing adoption in fisheries



The screenshot shows the GitHub repository page for 'nwfsc-assess / geostatistical_delta-GLMM'. The repository is described as a 'Tool for geostatistical analysis of survey data, for use when estimating an index of abundance'. It has 145 commits, 1 branch, 8 releases, and 2 contributors. The current branch is 'master'. The latest commit by James-Thorson is 'fixed bug in mean_D_tI computation' from 11 hours ago. The repository contains folders for 'R', 'data', and 'examples'.

Q1: Which spatial model is best?

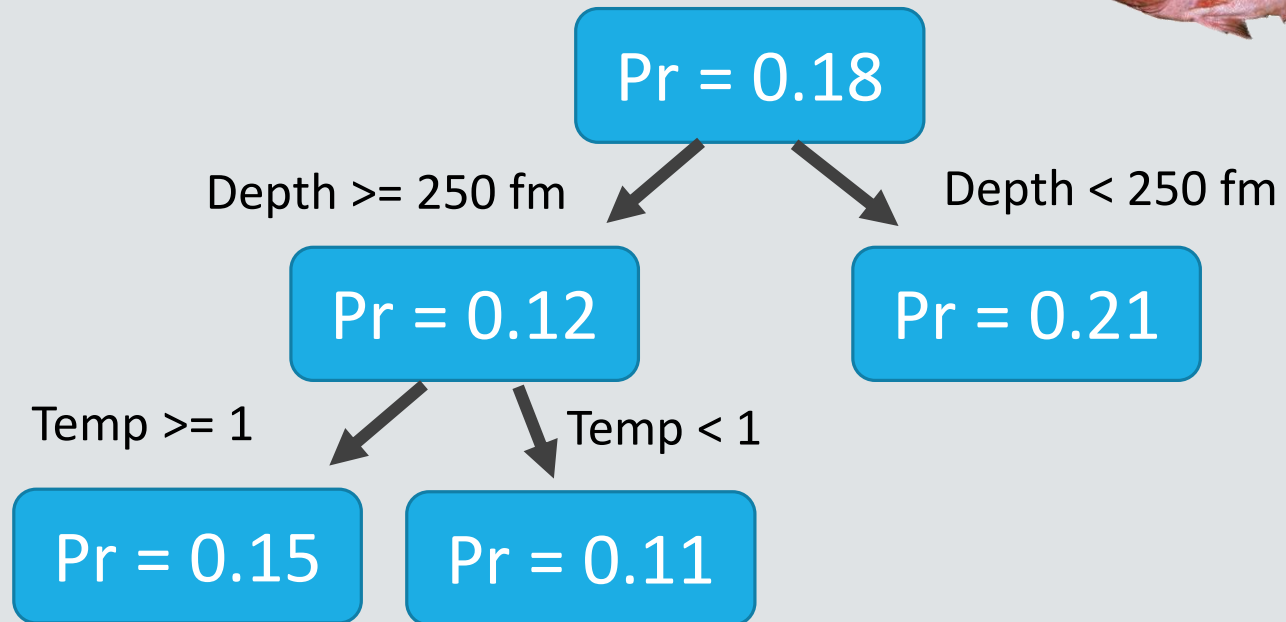
GLM

GAM

GMRF

RF

Random Forest



Q1: Which spatial model is best?

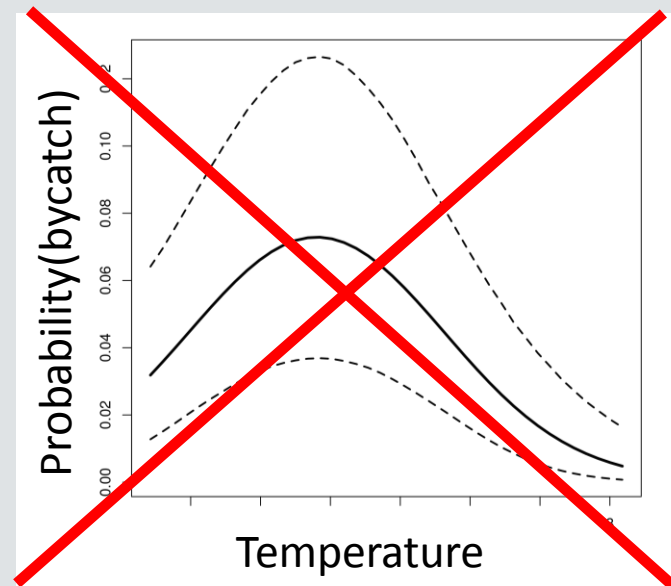
GLM

GAM

GMRF

RF

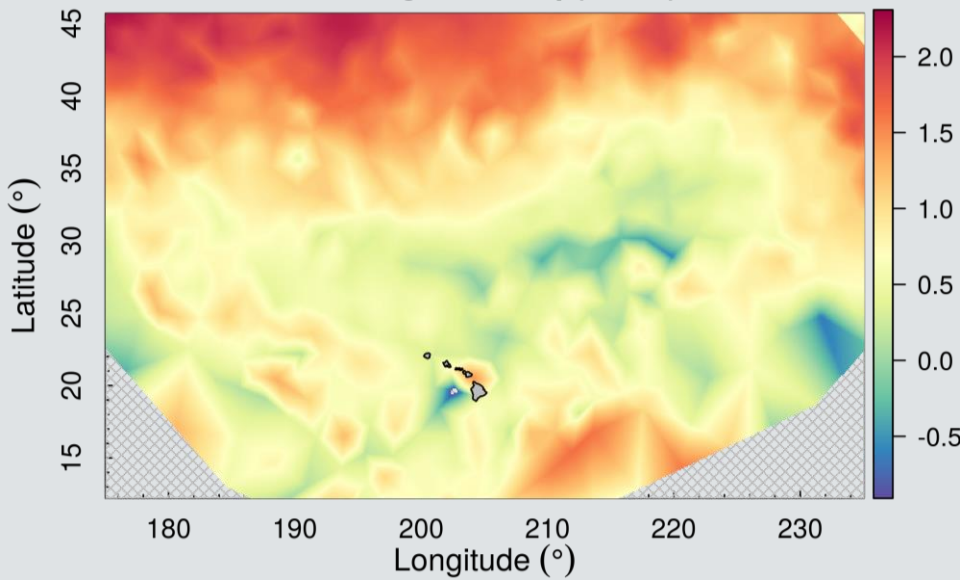
- Machine learning, designed for prediction
- “Black box”
 - Predictor-bycatch relationships not modeled
 - **No spatial field (add LAT, LON)**



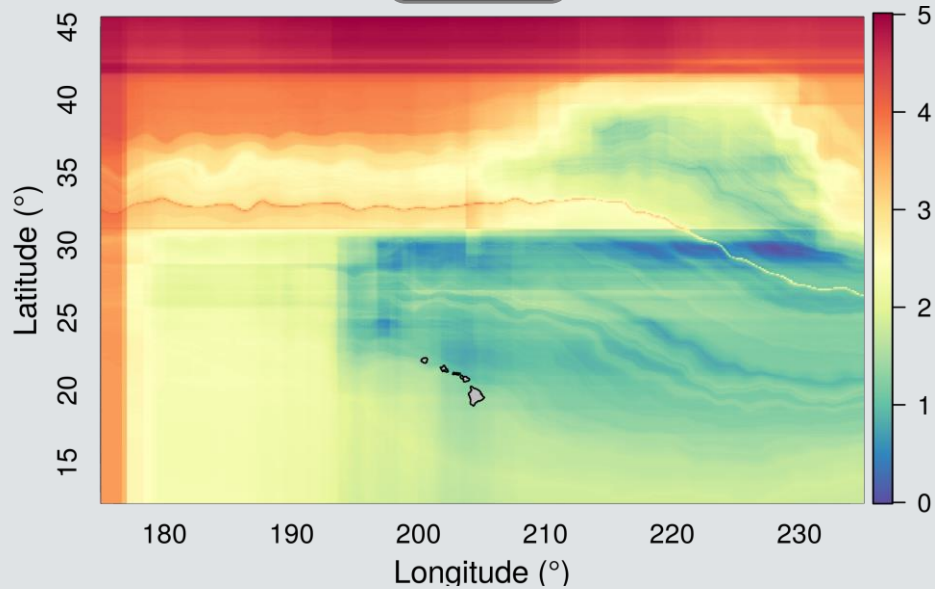
Bycatch risk maps



GMRF



RF



Results

Binomial

Generally:

GLM

<

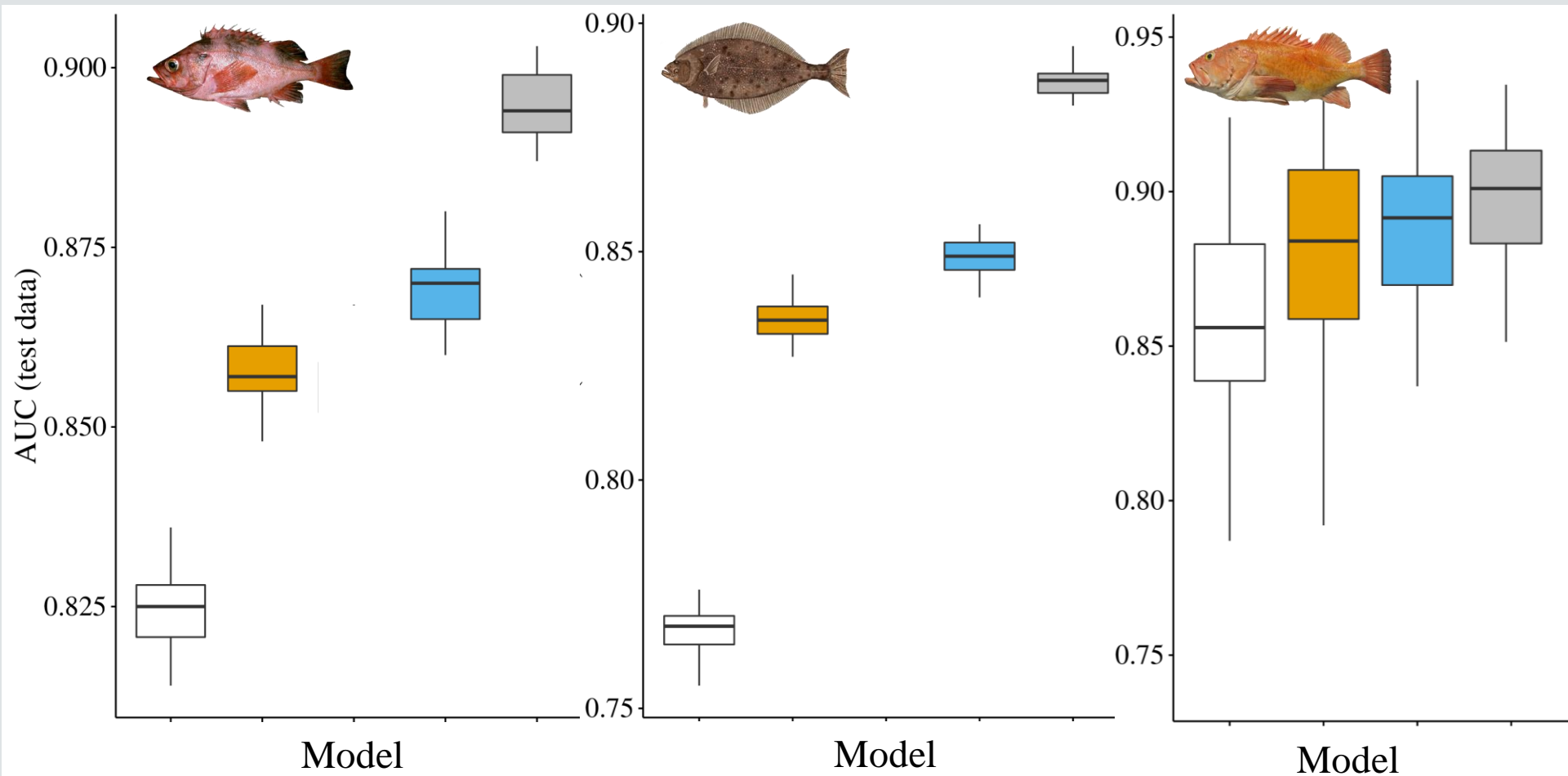
GAM

<

GMRF

<

RF



Results

Binomial

Generally:

GLM

<

GAM

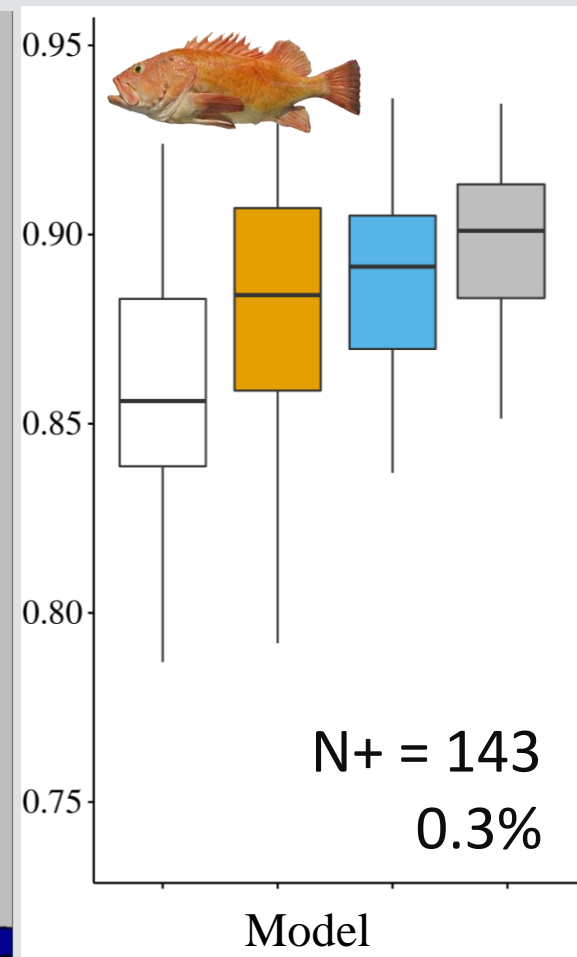
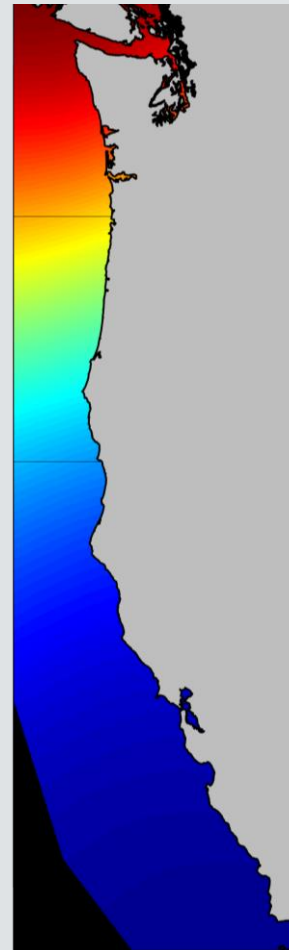
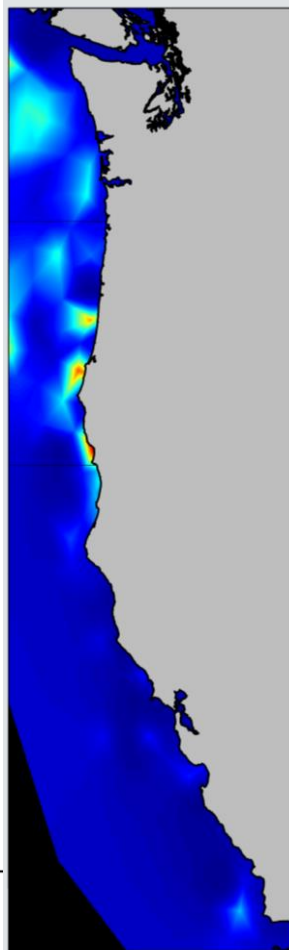
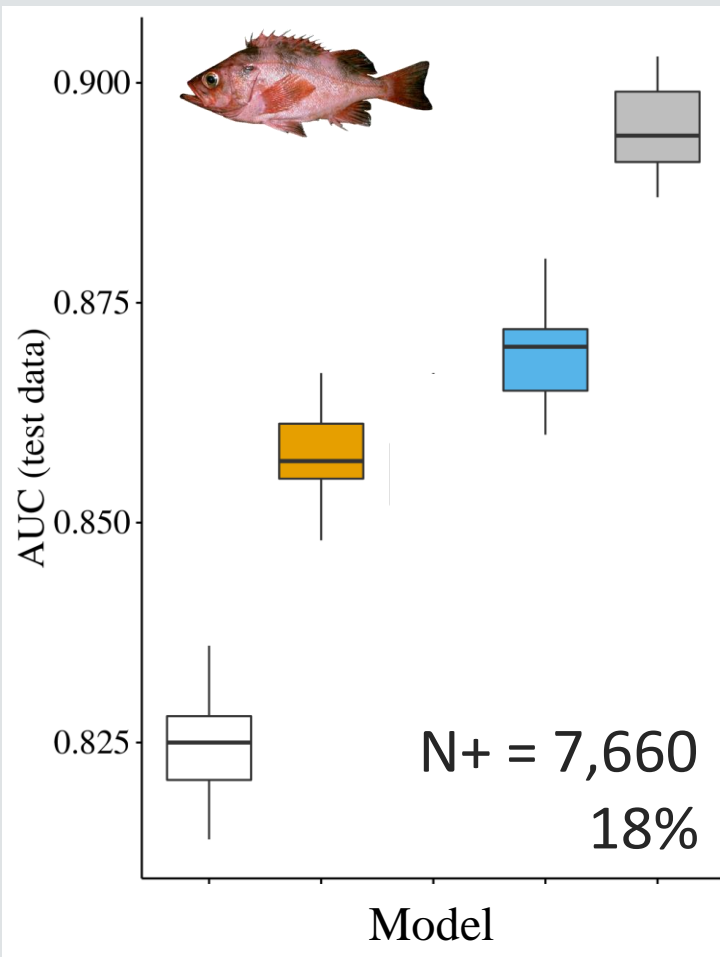
<

GMRF

<

RF

Less clear for rarer species



Results

Binomial

Generally:

GLM

<

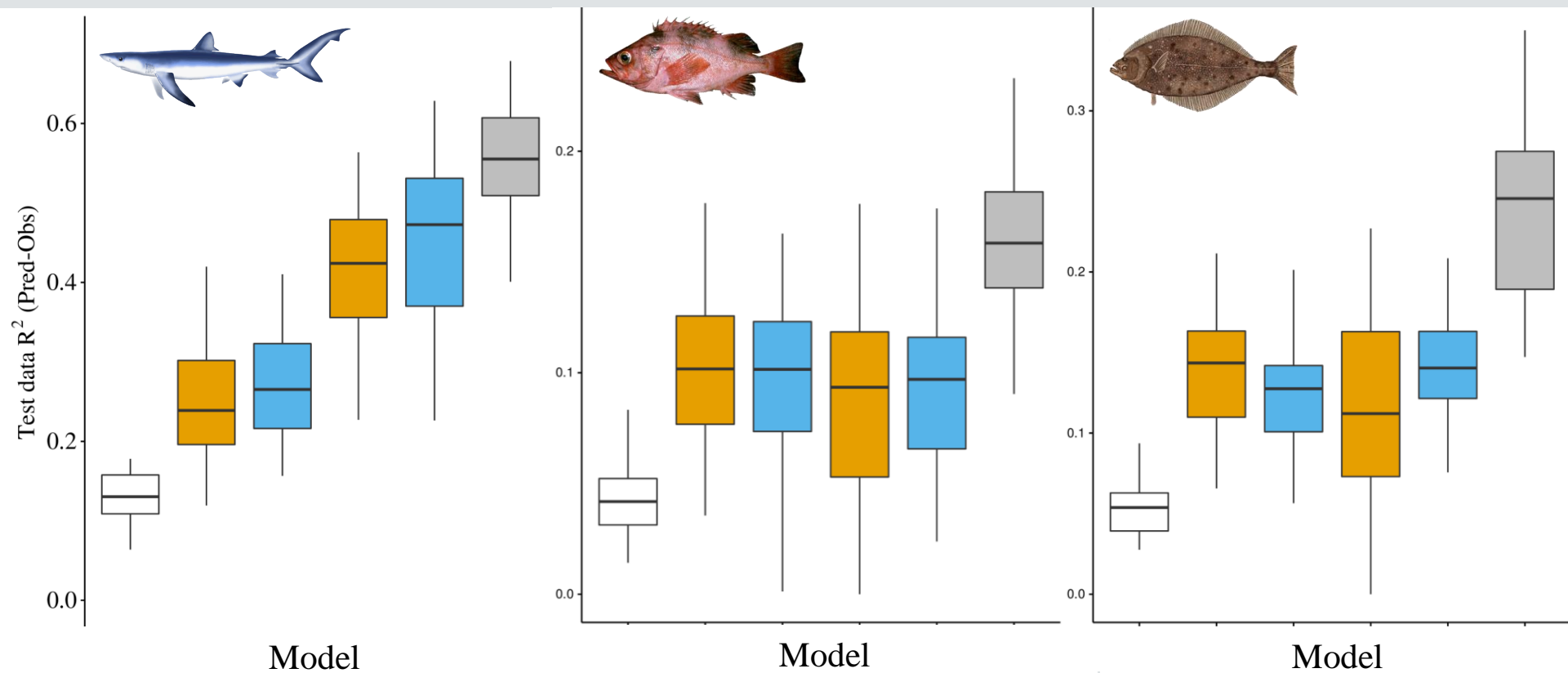
GAM

<

GMRF

<

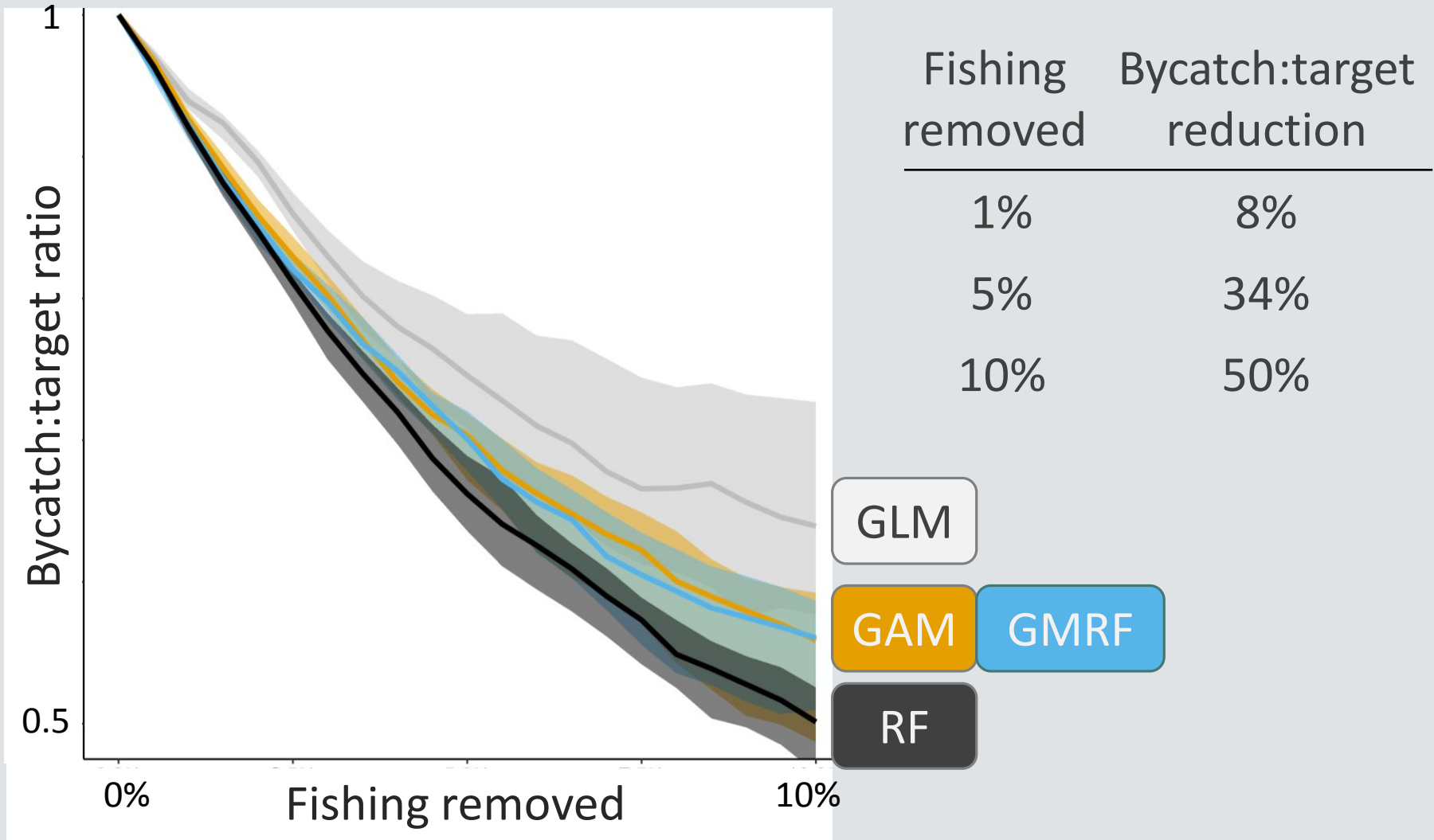
RF



Results

Positive

Q3: How much bycatch can they prevent?



Results

Conclusions

Q1: Which spatial model best predicts bycatch?



Q2: Does the answer depend on species?

No, **RF** had consistent advantage

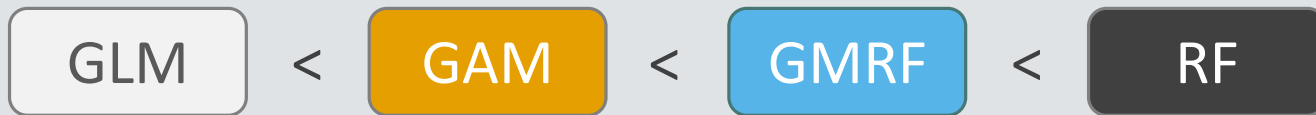
(larger for species with higher bycatch rates)

Q3: How much bycatch can they prevent?

Enough to consider using them in management

Discussion

If the goal is purely *prediction*:

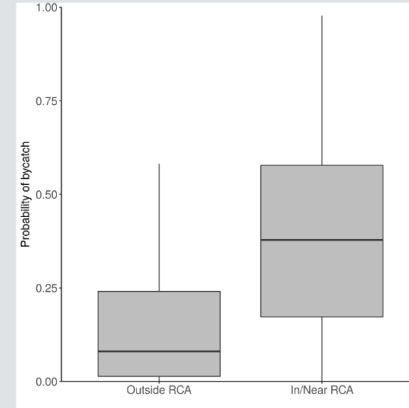
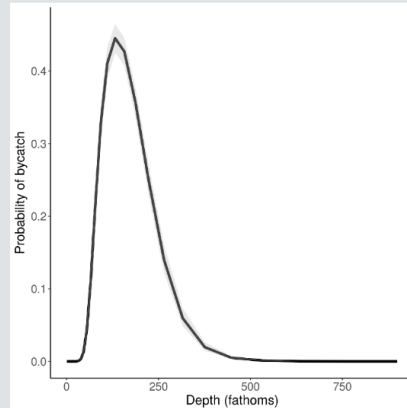
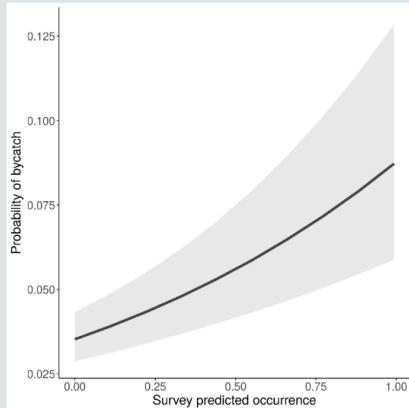


...but if we care about *inference on processes* affecting bycatch:

Covariate effects



GMRF

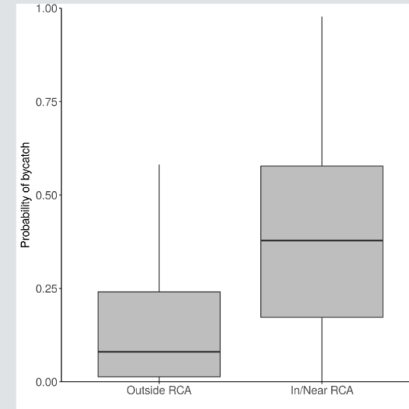
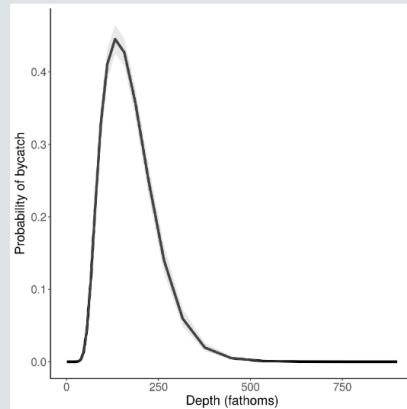
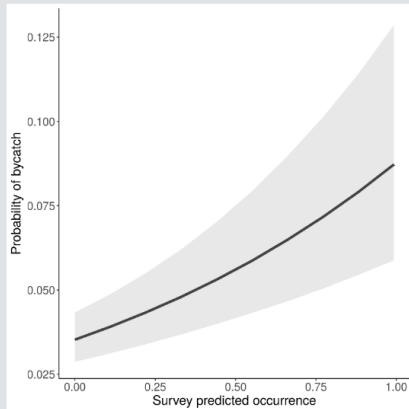


Covariate effects

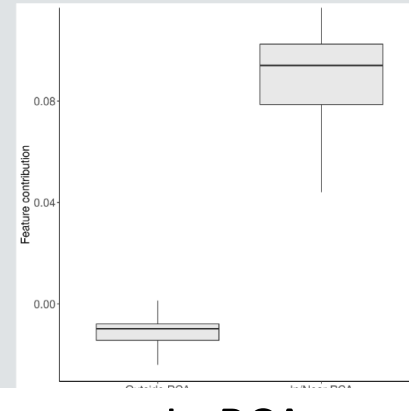
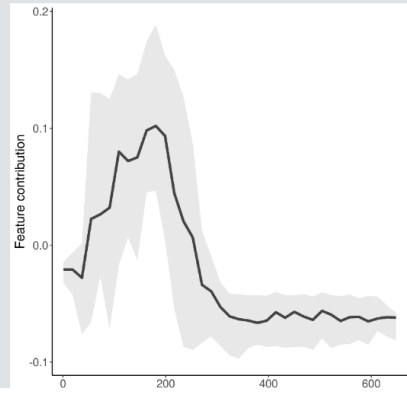
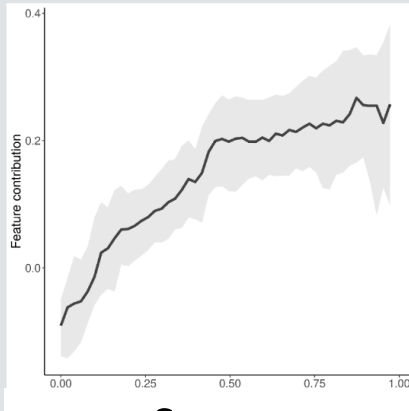
Are random forests really “black boxes”?



GMRF



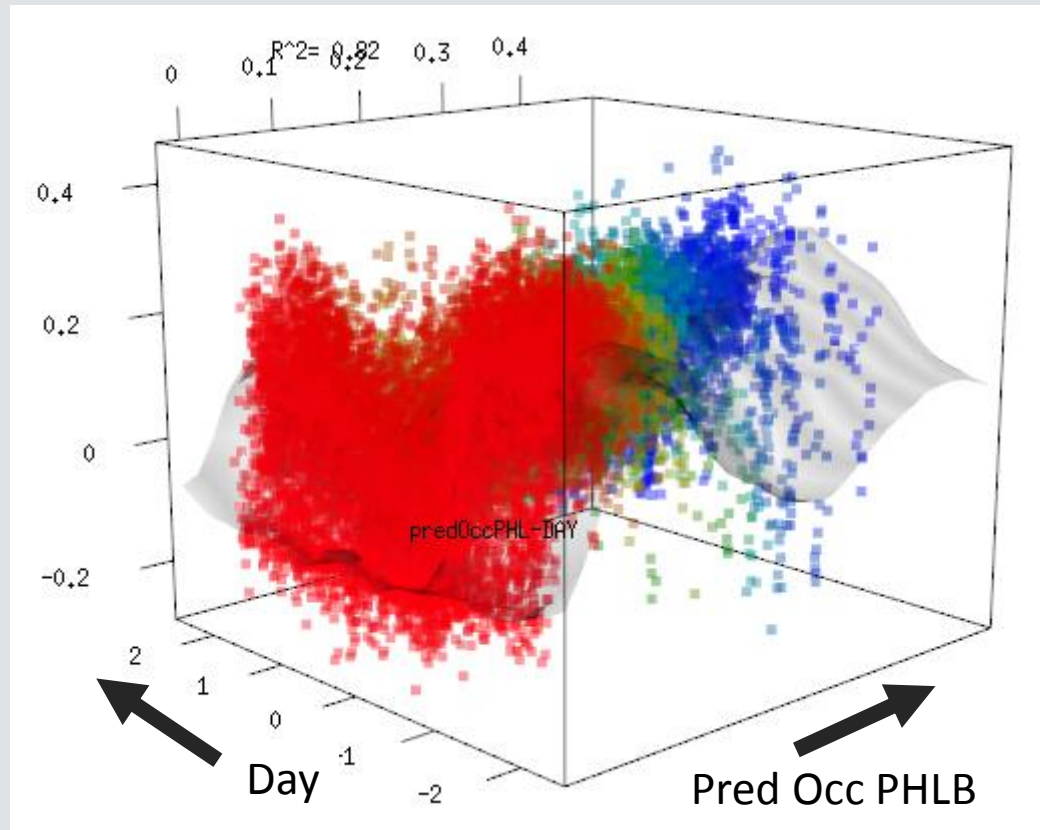
RF



Can random forests do *better*?

RF

Identifying covariate interactions



Discussion

Thank you!

SIO

- Brice Semmens

SWFSC

- Tomo Eguchi

NWFSC

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

PIFSC

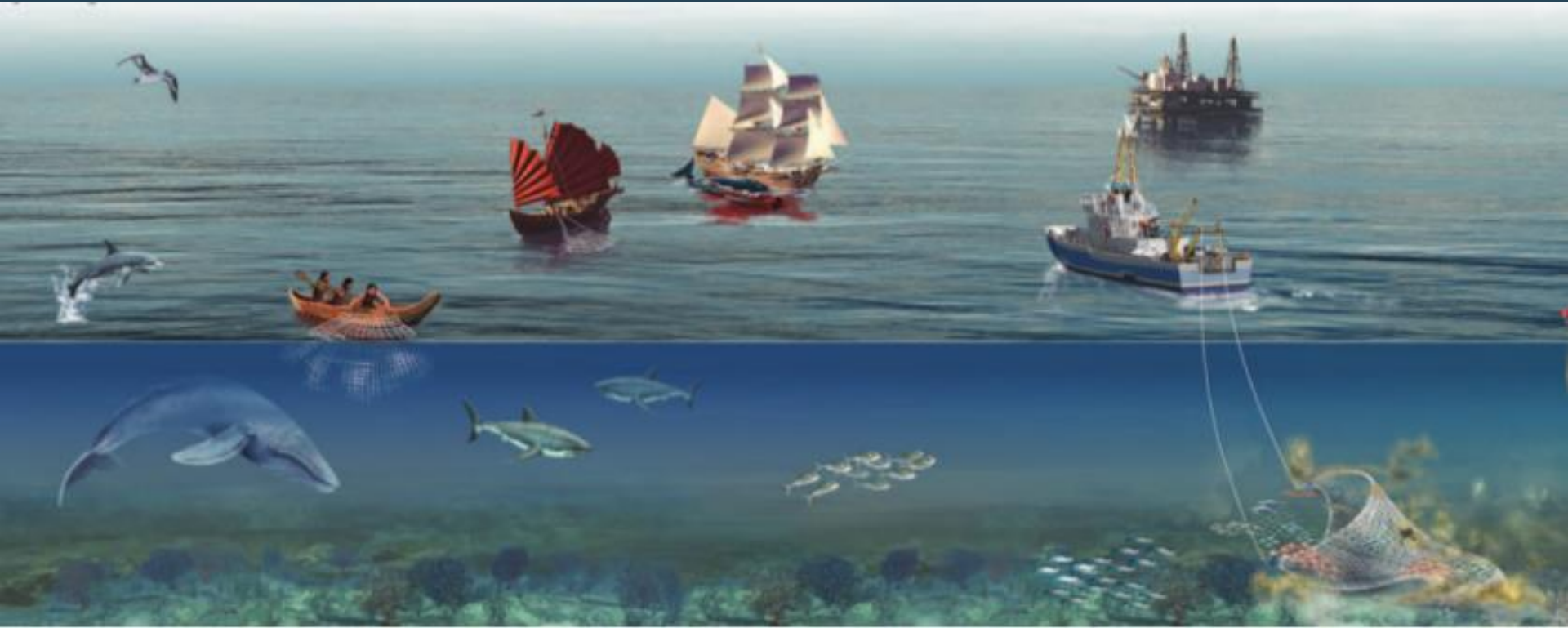
- Hawaii Longline Observer Program (Eric Forney)





Research opportunities in applied
math/statistics and fisheries science

We can easily harvest too many fish

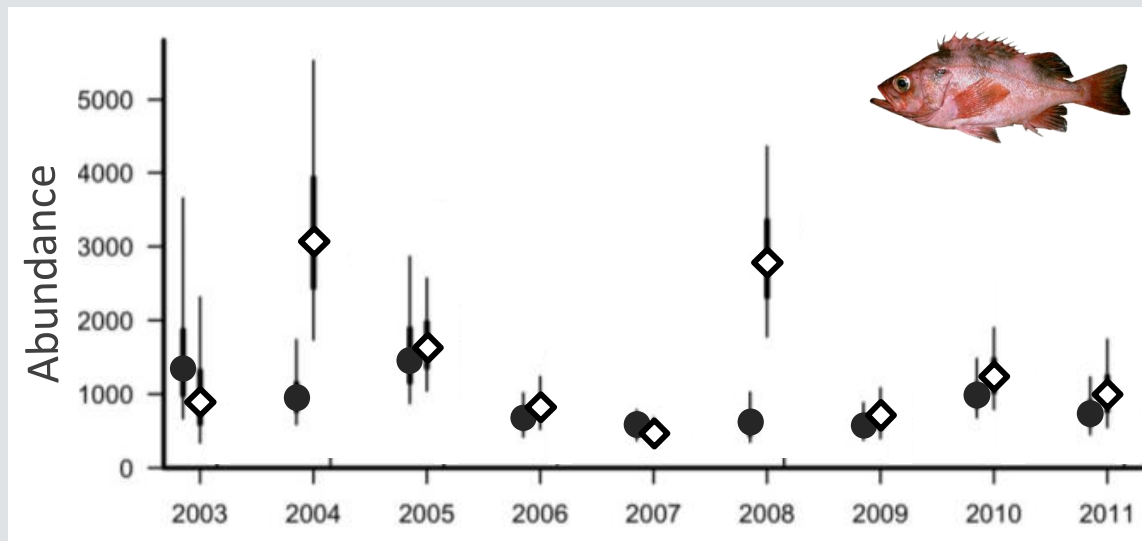


We use models in management

1. Sustainable harvest → need to assess populations

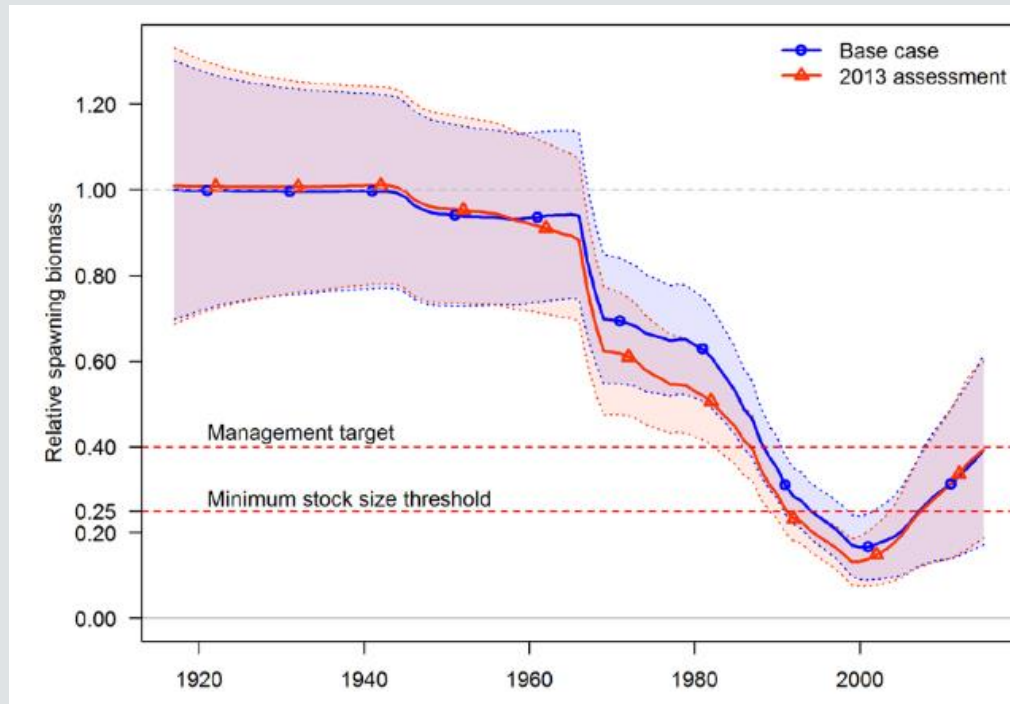
We use models in management

1. Sustainable harvest → need to assess populations
2. Primarily, *how many* and *where*



Build & test population models

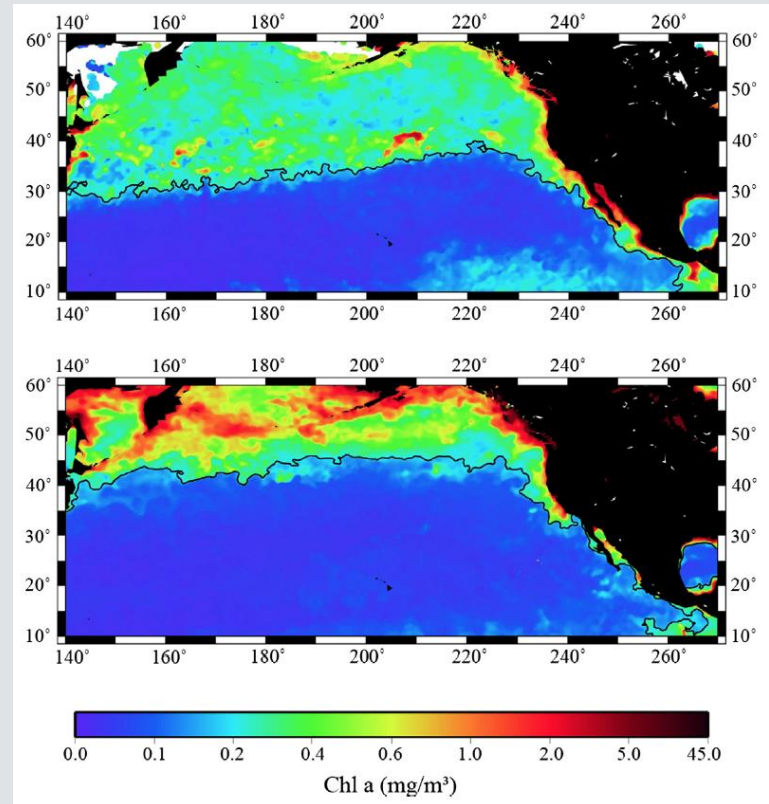
1. Stock assessment
2. Simulate alternative harvest strategies



Work with physics/climate modelers

What are the effects on fish of:

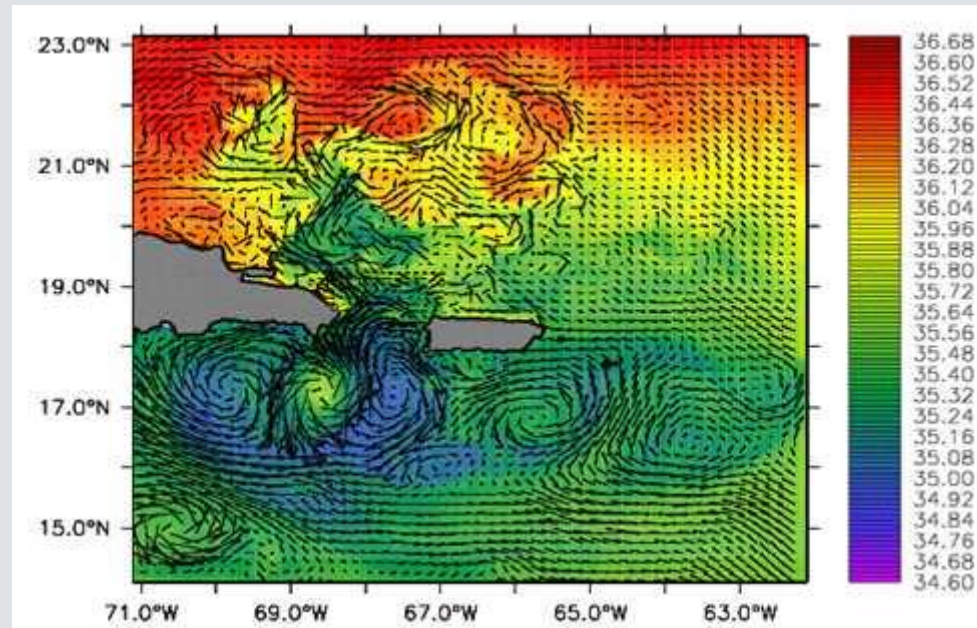
1. Ocean productivity?



Work with physics/climate modelers

What are the effects on fish of:

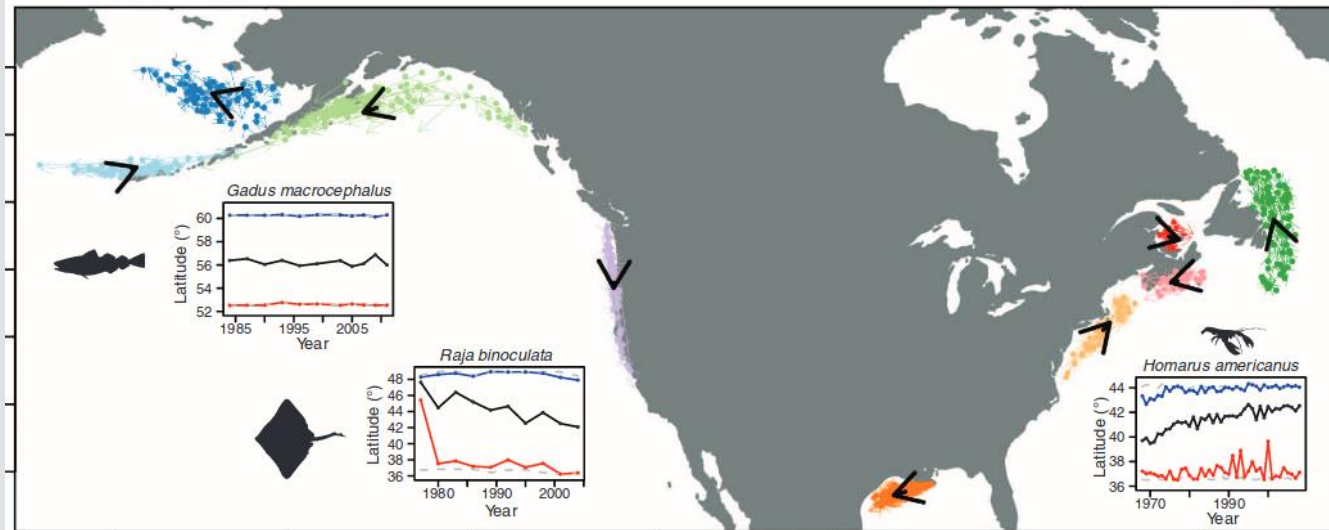
1. Ocean productivity?
2. Dispersal of eggs and larvae?



Work with physics/climate modelers

What are the effects on fish of:

1. Ocean productivity?
2. Dispersal of eggs and larvae?
3. Range shifts?



How to gauge model performance?

Goal: prediction

5-fold cross validation repeated 10x

Binomial

ROC curve (AUC)

Positive

RMSE

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)
- + Much quicker to set up and run (~2 min vs. 5-15 hours)
- + Not just a “black box”?

GMRF

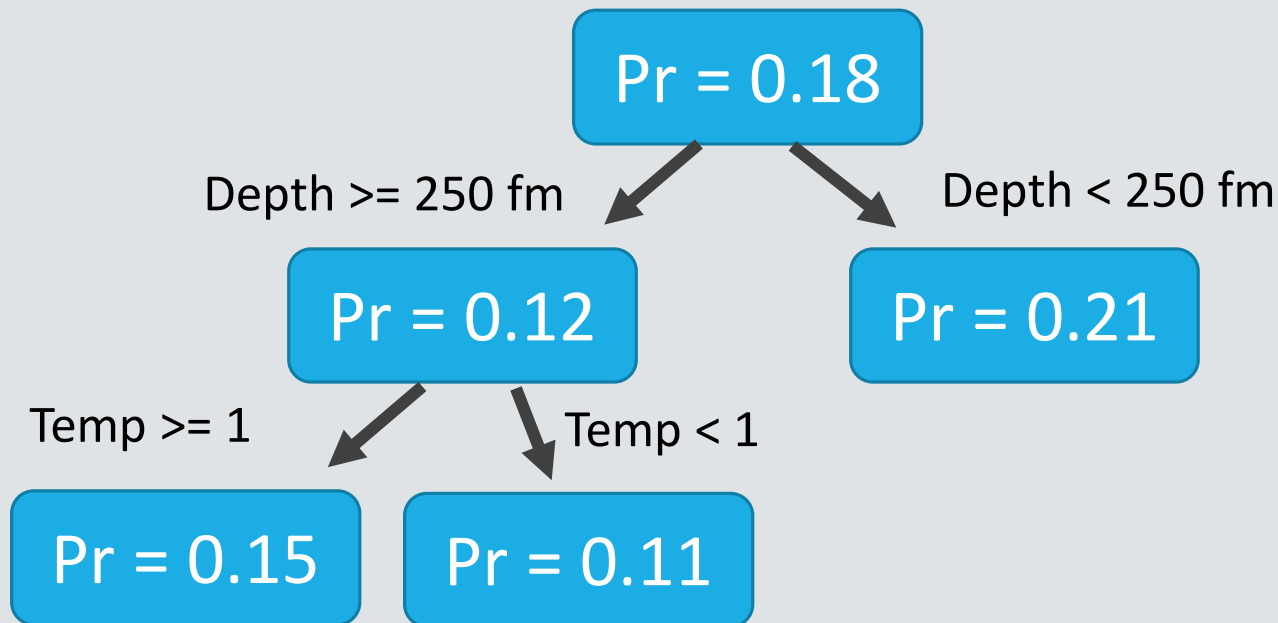
- + Statistical inference, marginal posteriors for covariate effects
- + Ability to include observation error

Covariate effects

RF



What is a “feature contribution”??

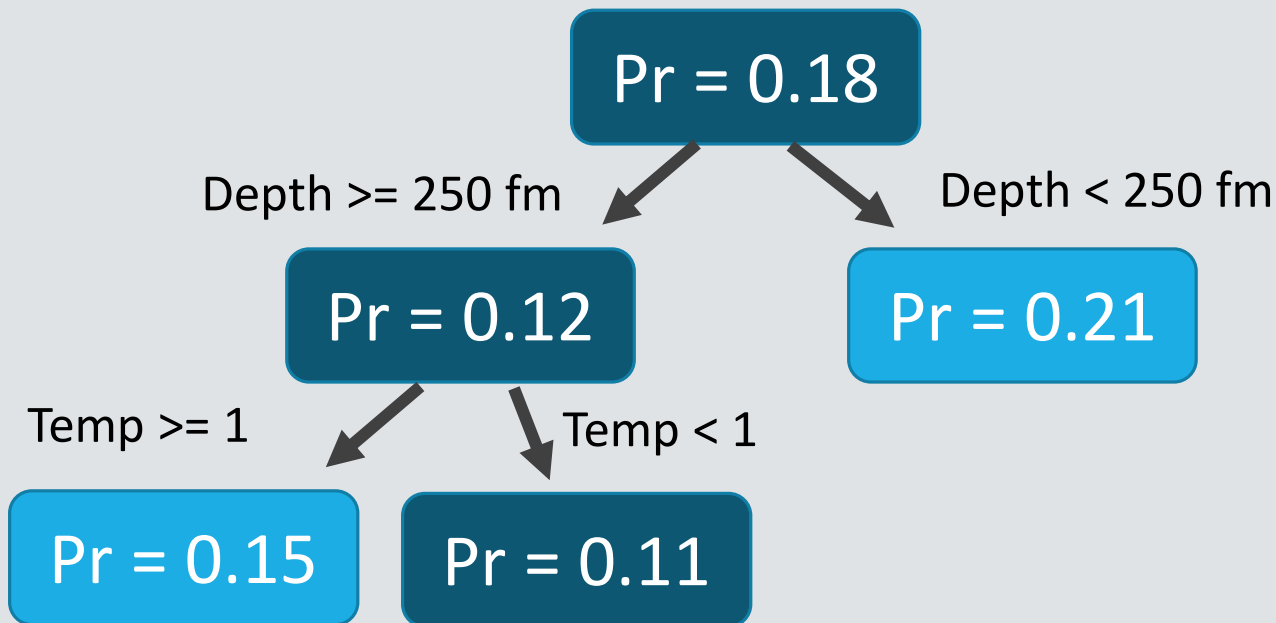


Covariate effects

RF



What is a “feature contribution”??



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