

NLDAS-3 Telecon @ 2015.07.15



# Statistical Precipitation Downscaling Using Random Forests

- *Synthetic experiments over Southeast United States* -

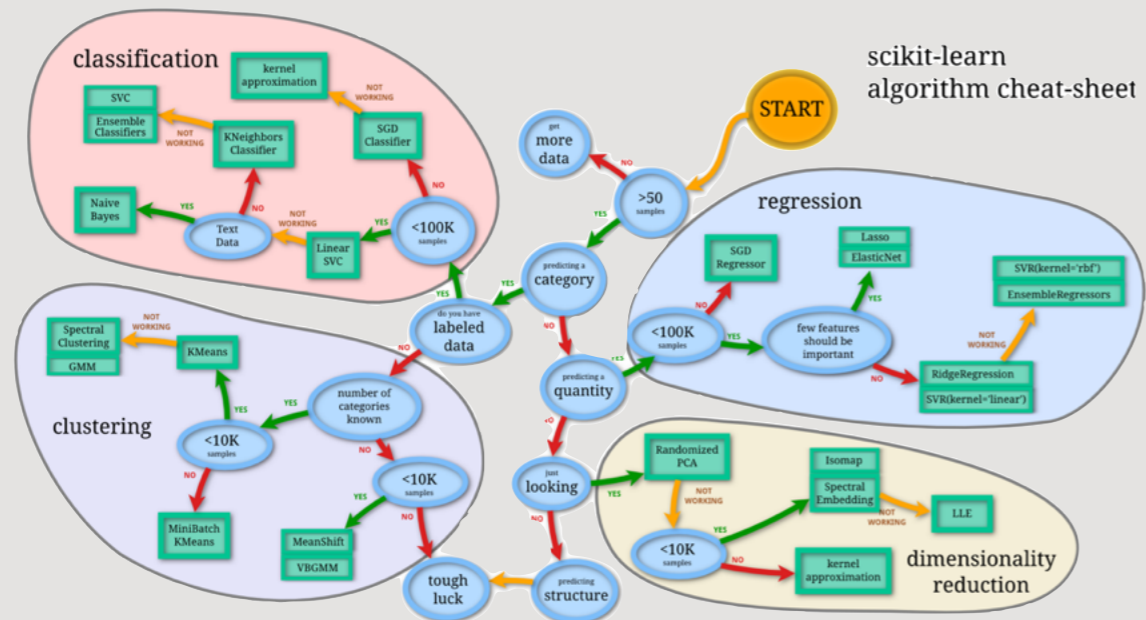
**Xiaogang He<sup>1</sup>, Nathaniel W. Chaney<sup>1,2</sup>, Justin Sheffield<sup>1</sup>,  
Ming Pan<sup>1</sup>, Eric F. Wood<sup>1</sup>**

<sup>1</sup> *Department of Civil and Environmental Engineering, Princeton University*

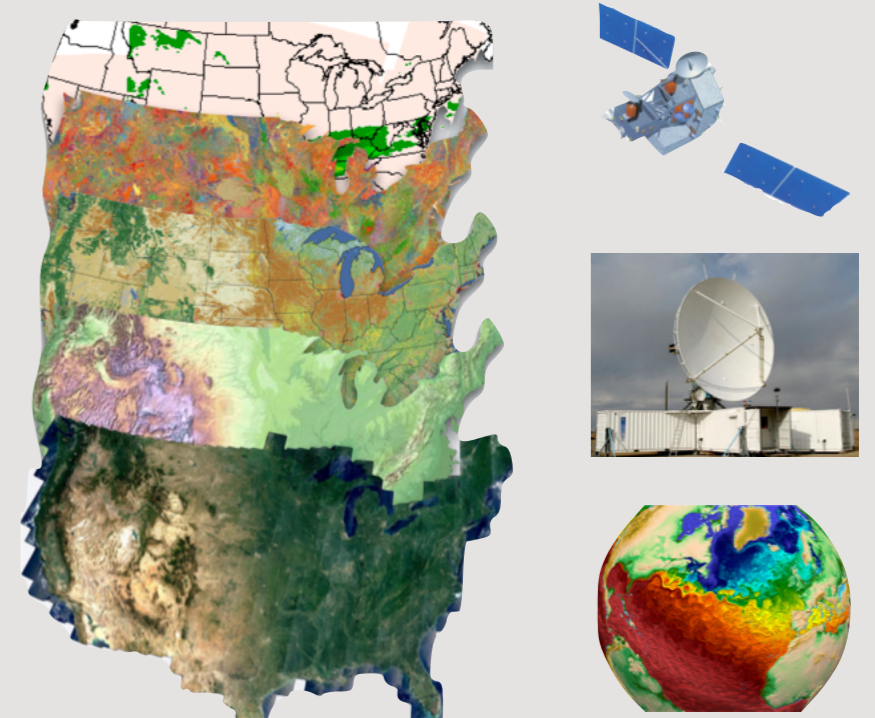
<sup>2</sup> *Program of Atmospheric and Oceanic Sciences, Princeton University*



## Machine learning approach



## Big data



Explore and construct complicated relationships

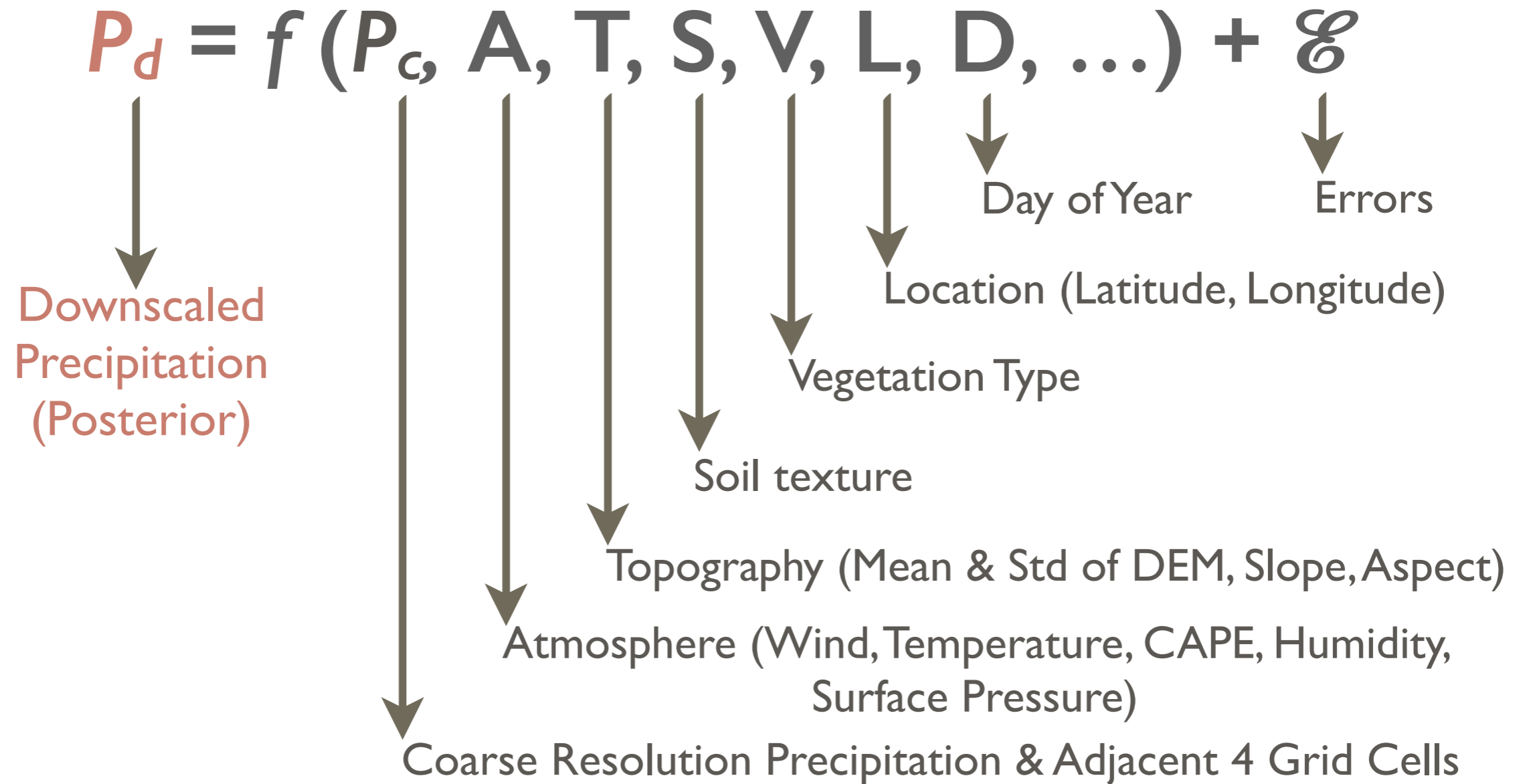


High performance computing



Precipitation spatial-temporal structures??

**Hyperresolution**  
precipitation downscaling??



## Single learner



**VS**

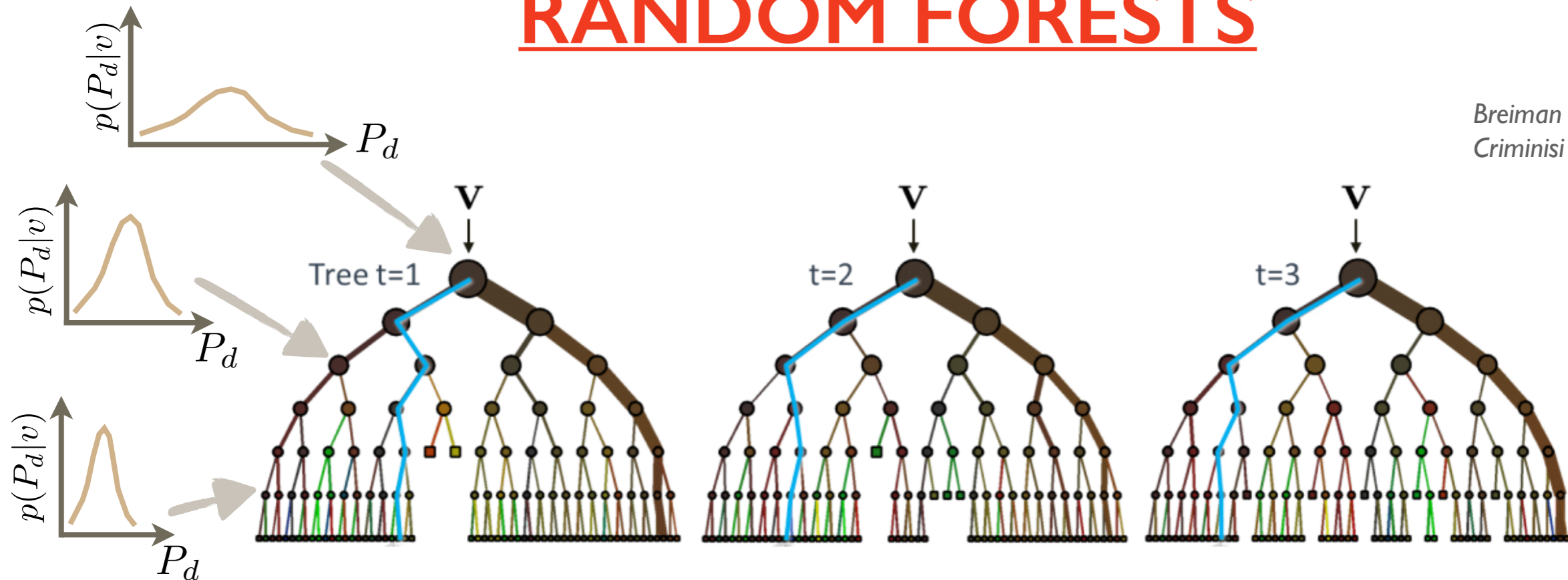
## Group of learners



Ensemble methods → Reduce overfitting

# RANDOM FORESTS

Breiman (2001)  
Criminisi et al., (2011)



Regression forest posterior: 
$$p(P_d|\mathbf{v}) = \frac{1}{T} \sum_{t=1}^T p_t(P_d|\mathbf{v})$$

$\mathbf{v}$ : covariates

$P_d$ : predictand (downscaled prec)

$p_t(P_d|\mathbf{v})$ : individual tree posterior

# Step 1: Generate synthetic covariates

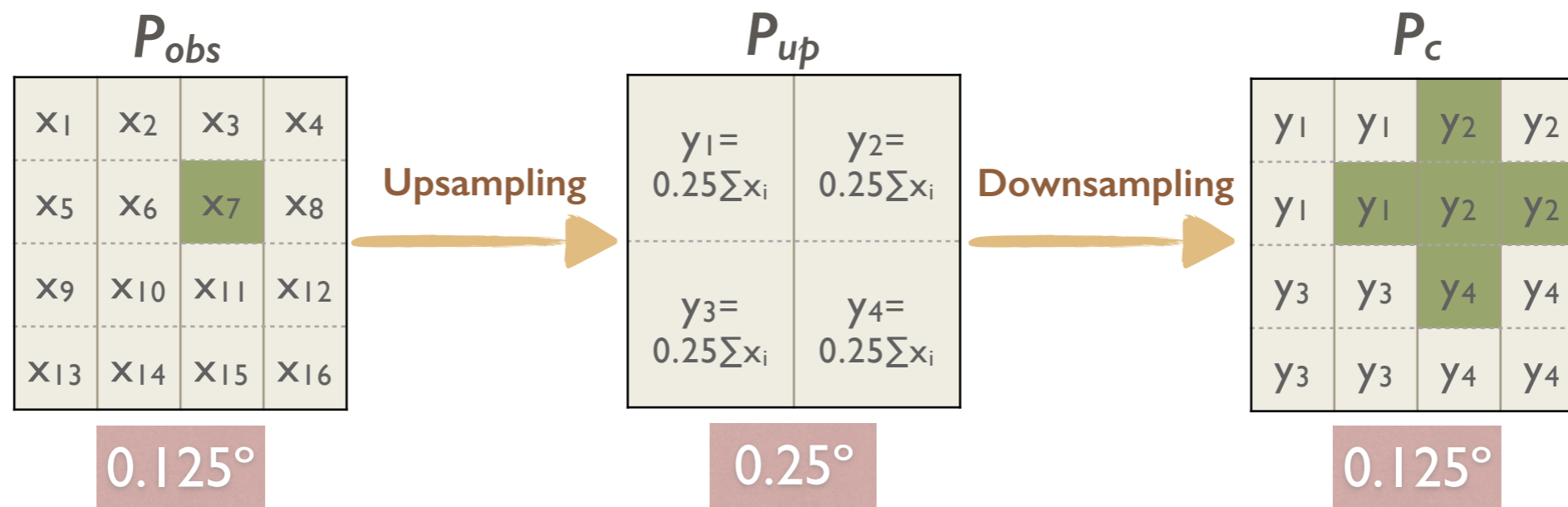
# Synthetic Experiment Design

**Data:** NLDAS2  
 North American Land Data Assimilation System 2

**Period:** 2011.06.01—2011.08.31 (Hourly)

**Domain:** Region: Southeast United States  
Grid: 80×72 (0.125°)

EXP	$P_C$	A
$P_{0.25}A_{0.125}$	0.25°	0.125°
$P_{0.25}A_{0.25}$	0.25°	0.25°
$P_{0.5}A_{0.125}$	0.5°	0.125°
$P_{0.5}A_{0.5}$	0.5°	0.5°
$P_1A_{0.125}$	1°	0.125°
$P_1A_1$	1°	1°



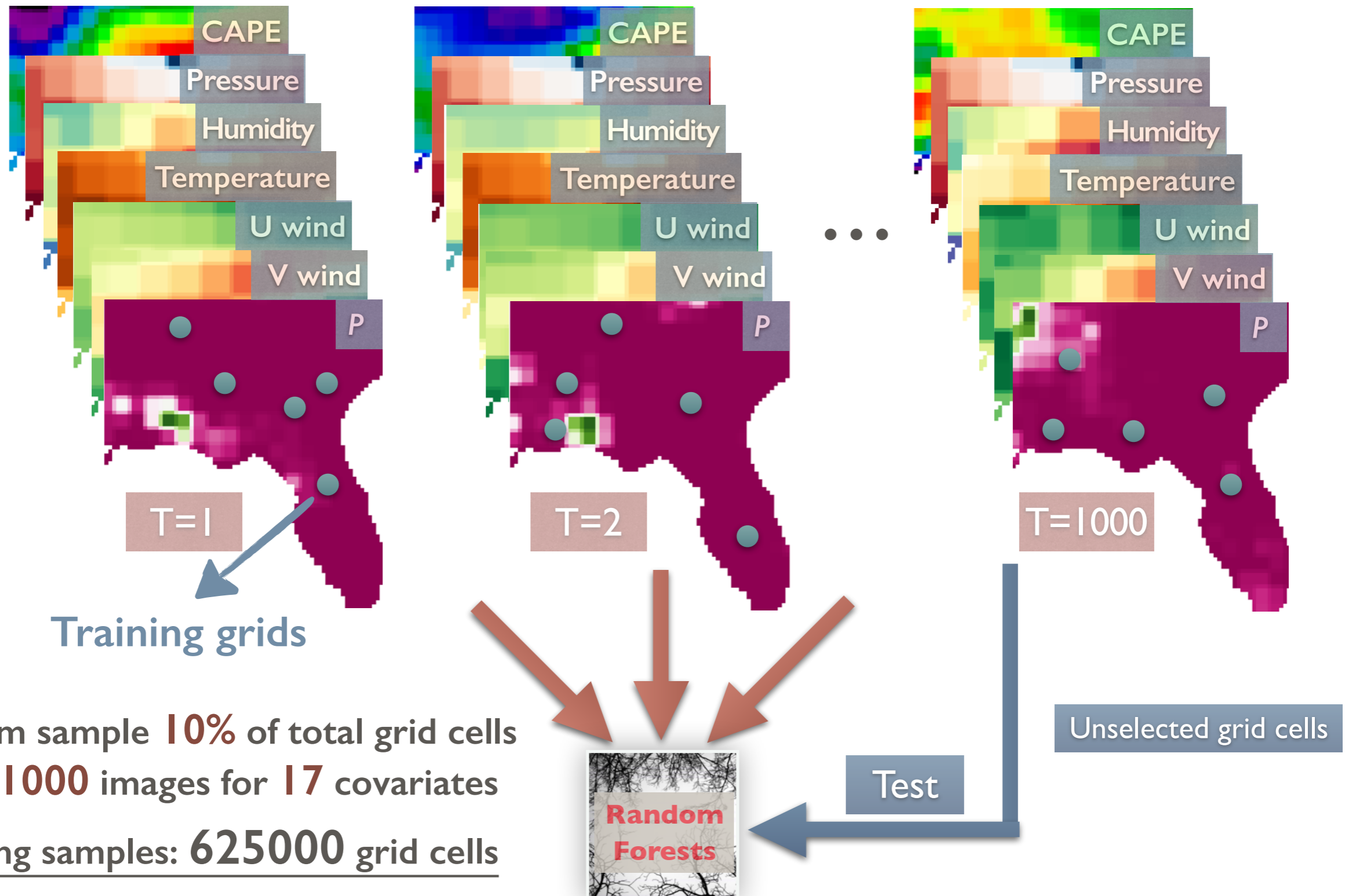
Adjacent grid cells are also considered  
 Same procedures for atmospheric covariates

# Step 2:

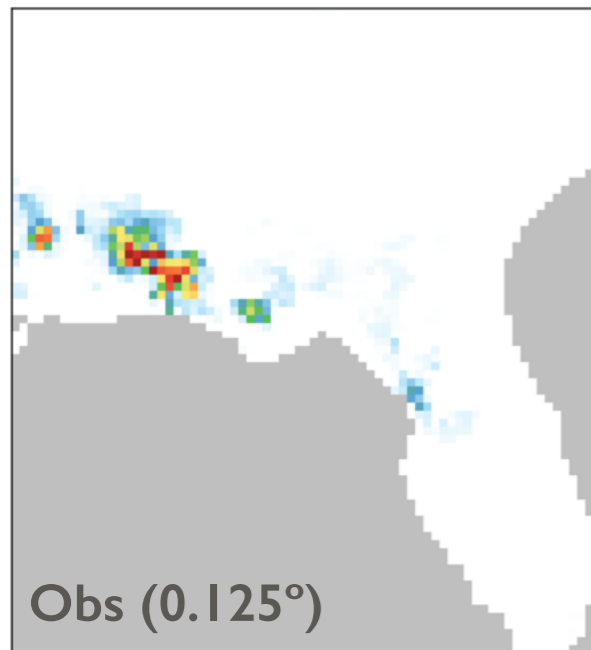
## Train and test Random Forests

# Synthetic Experiment Design

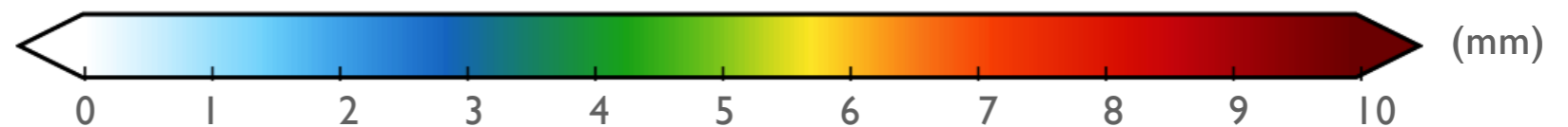
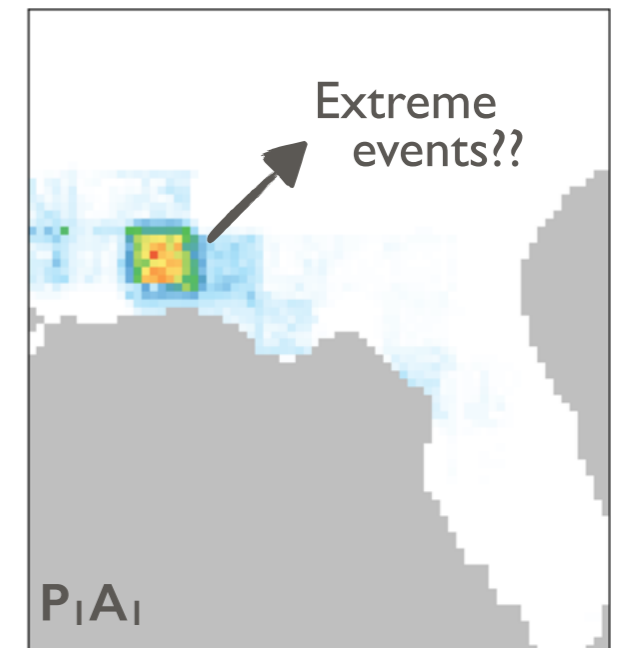
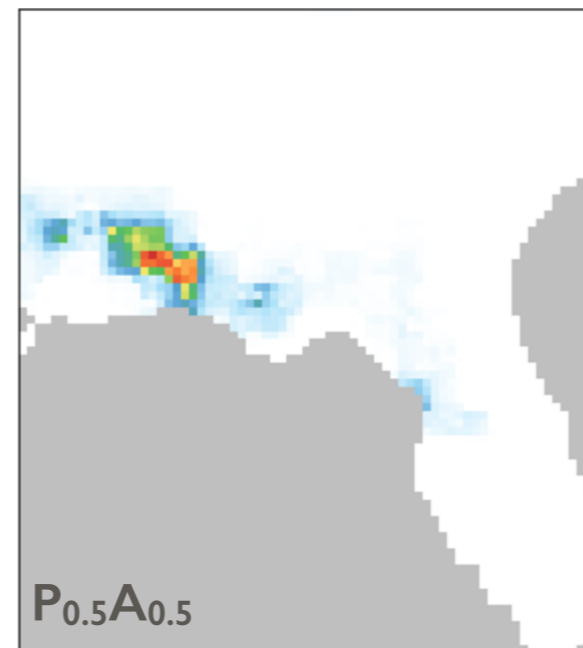
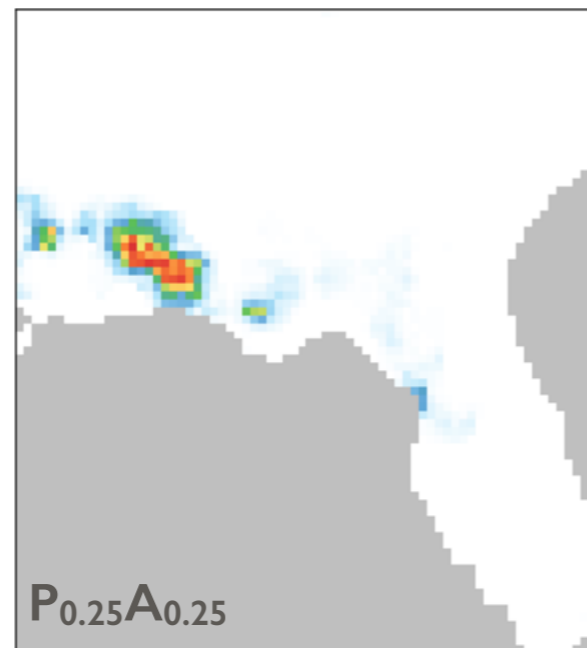
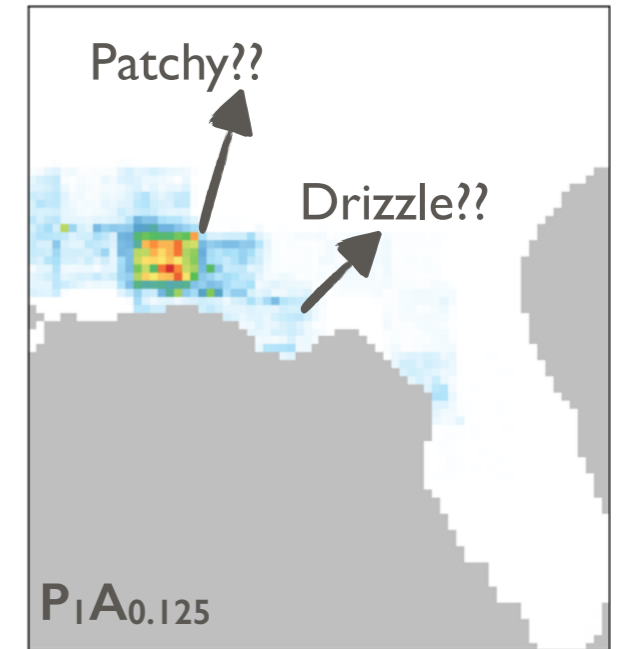
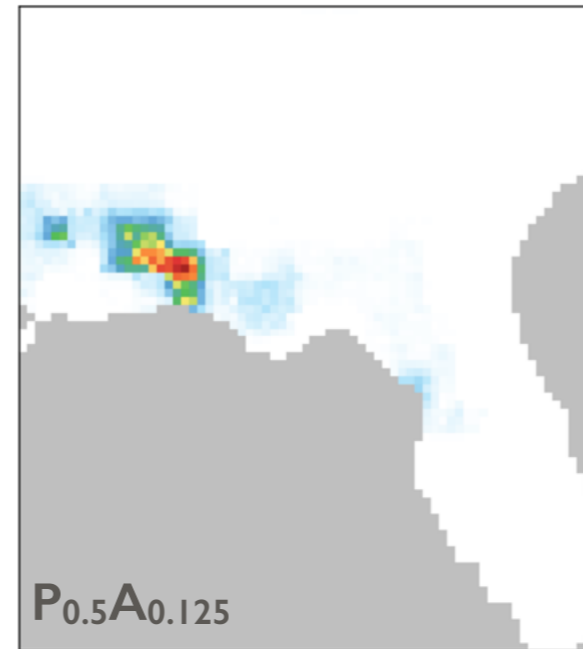
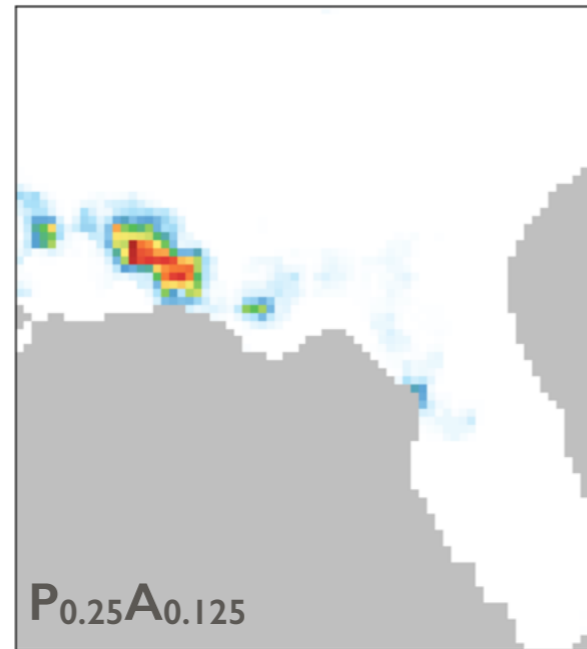
Coarse  $P$  + Dynamic/Fixed Covariates  $\Rightarrow$  Downscaled  $P$



## Scaling Experiments

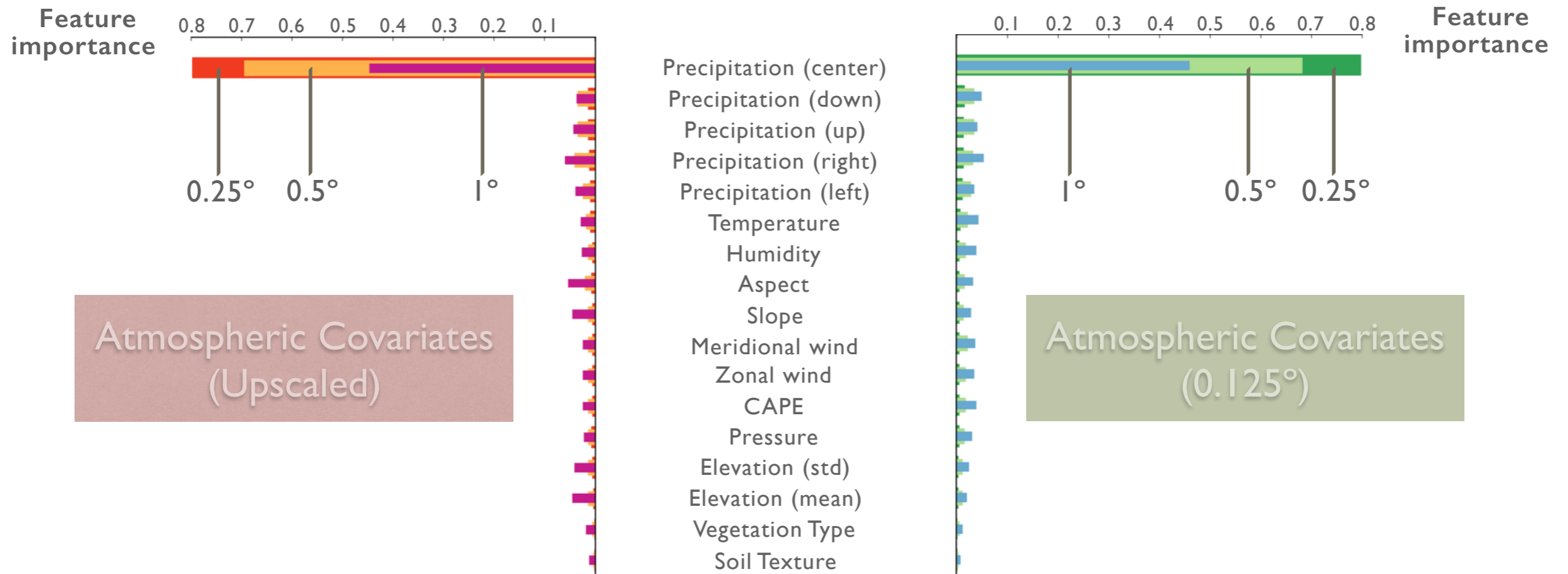


Example:  
2011.06.06.23:00



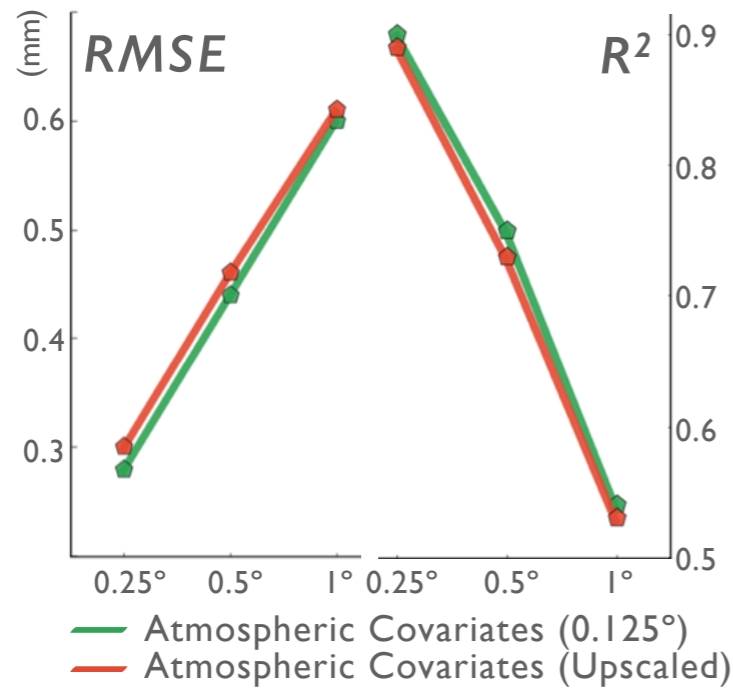


Feature importance measures the prediction strength for each covariate

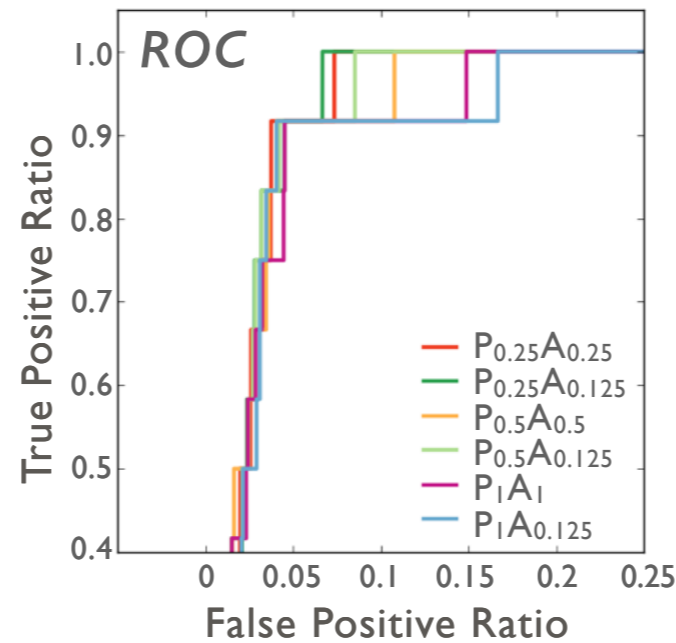


- ★ Central coarse precipitation dominates
- ★ Dynamic fields matter
- ★ Topography matters (e.g.,  $P|A|$ )

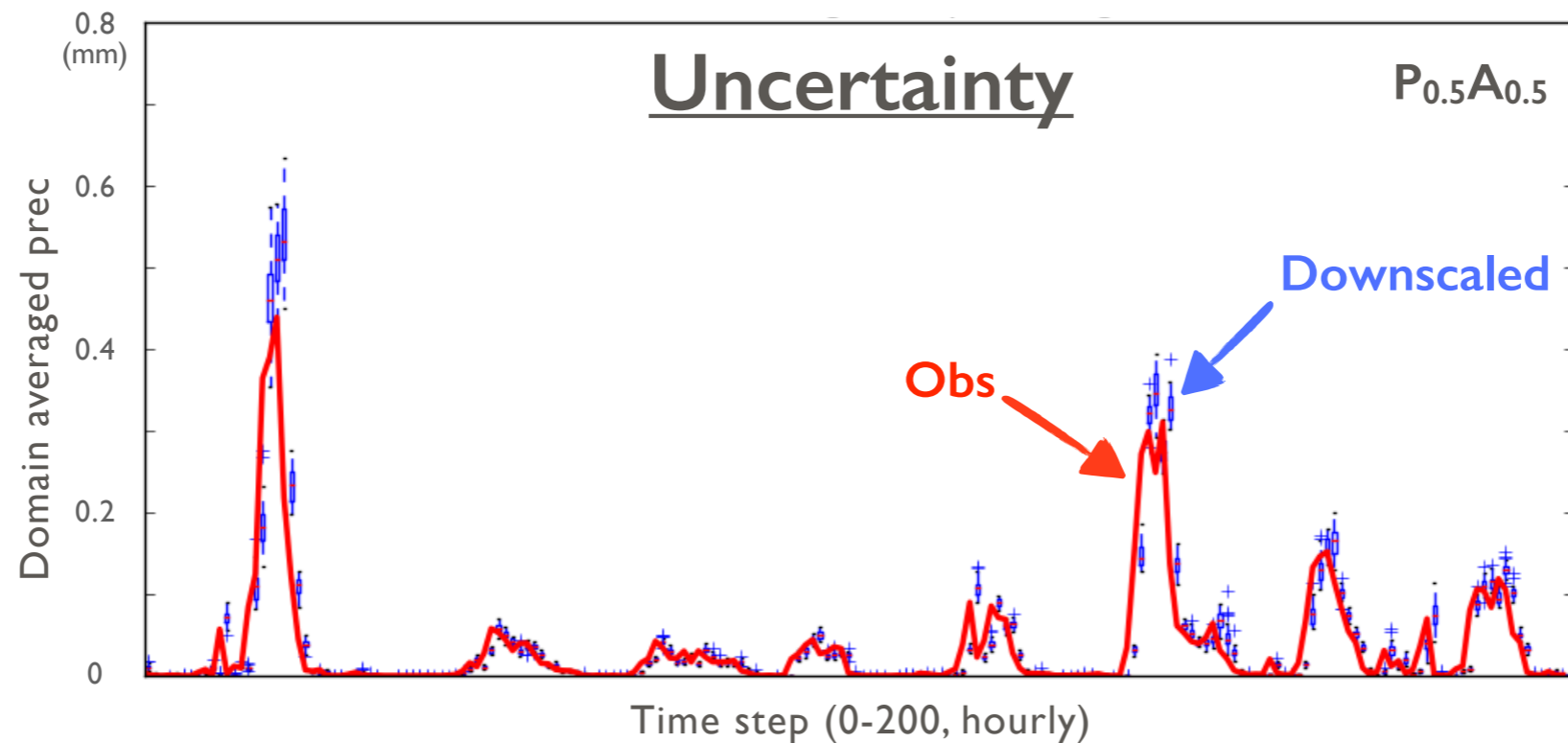
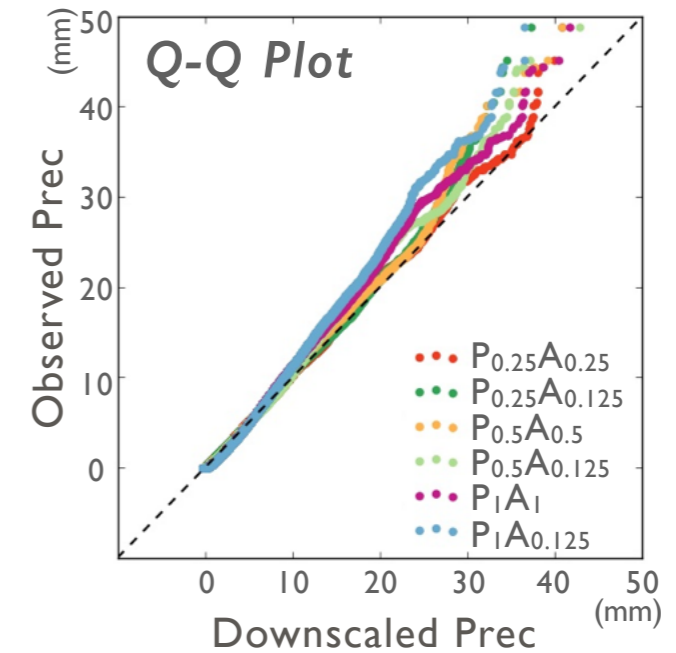
## Goodness-of-fit



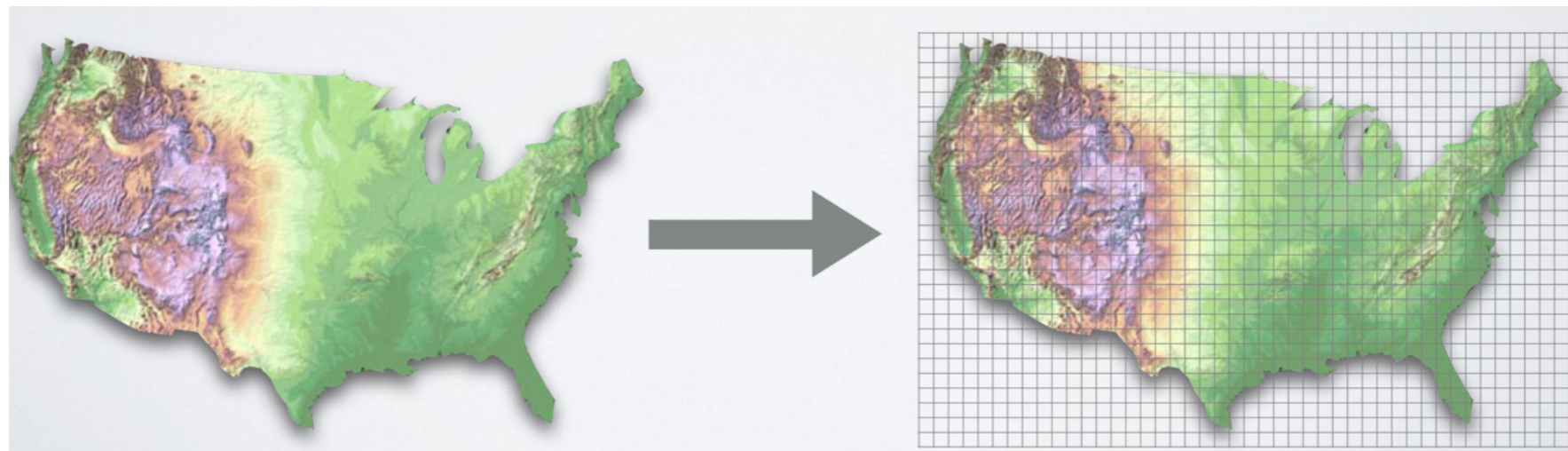
## Classification



## Distribution



- Synthetic experiments → Try other regions
  - Western US (California and mountainous regions)
  - Northeastern climate division
  - Central US
- Real experiments
  - Train StageIV, downscale satellite/reanalysis
- CONUS/Global scale → Moving window approach



*Thanks!!*

*Questions??*