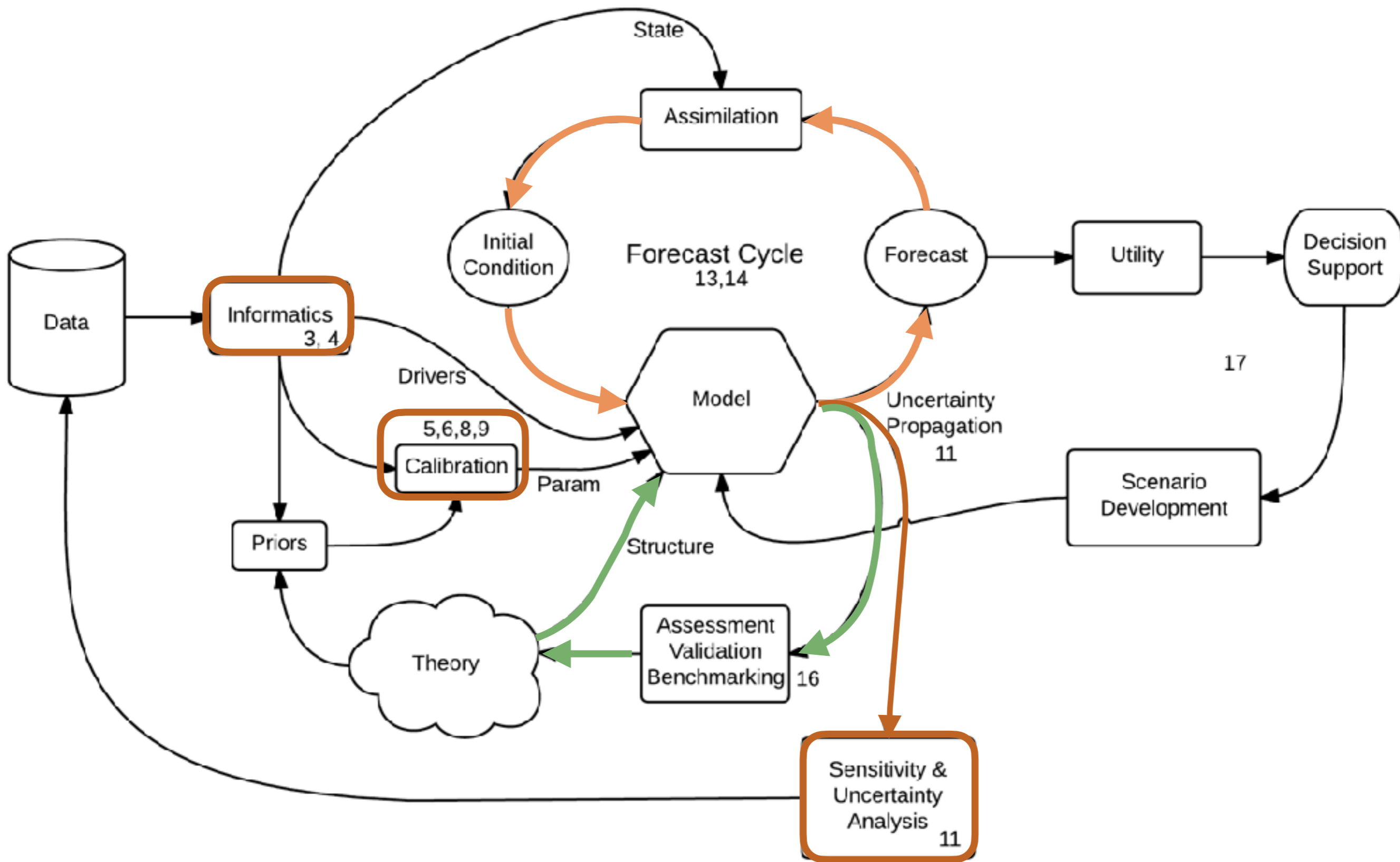# Assessing Model Performance

*Lesson 11*

- Range of values
- Units
- General pattern in time & space
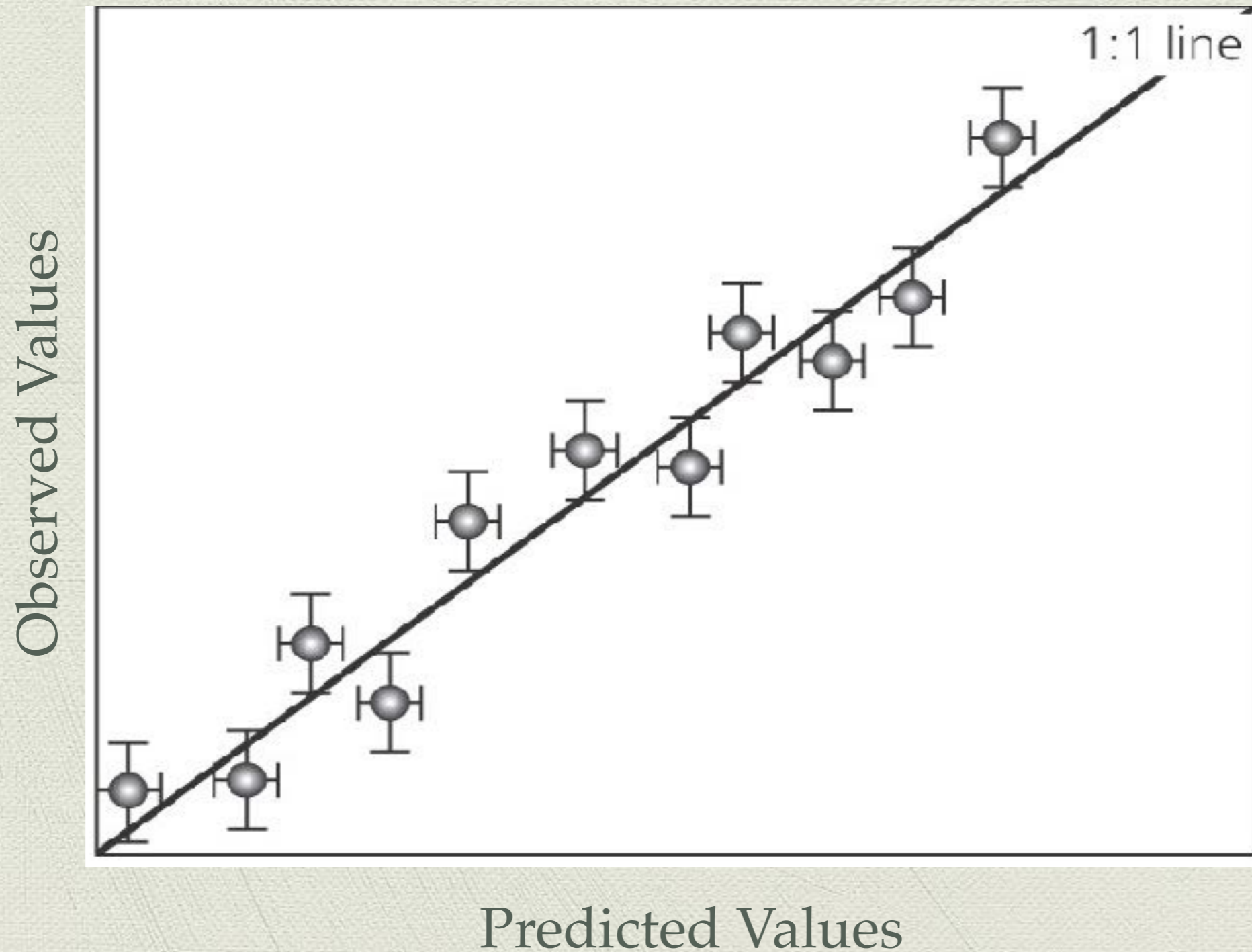
*Sanity check!*
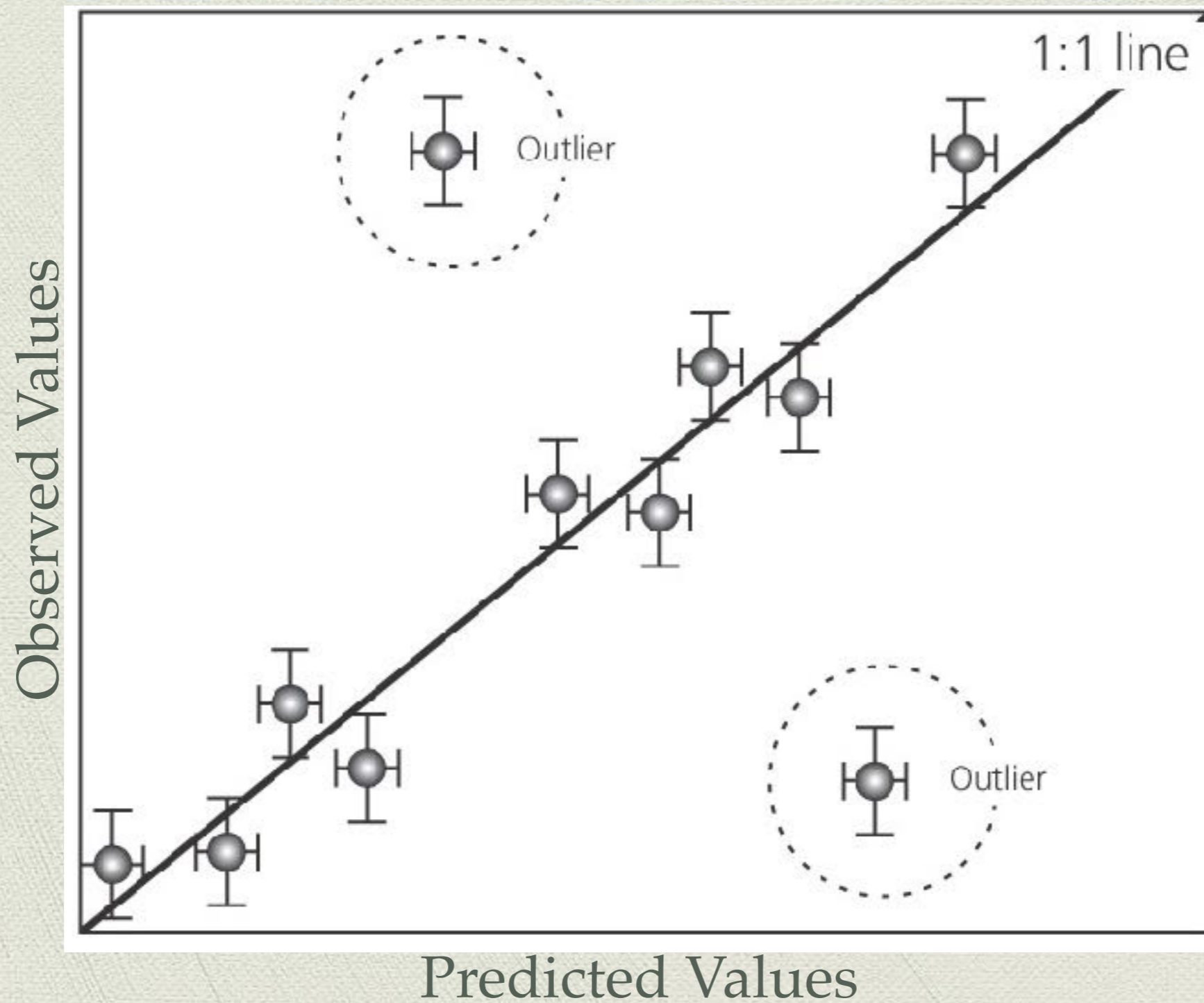
# Step 1: Is the model output reasonable?
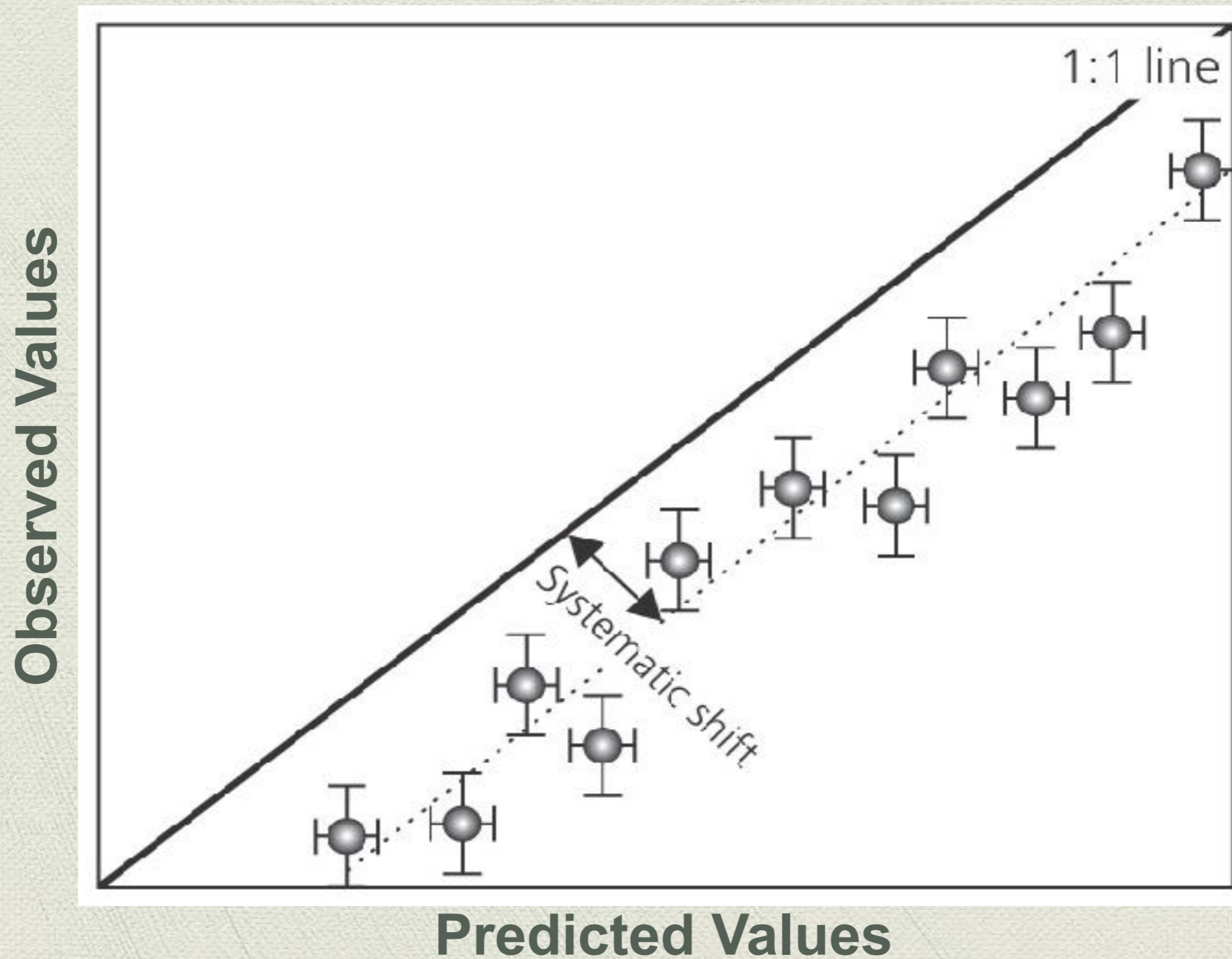
# Step 2: Graphical comparisons to data
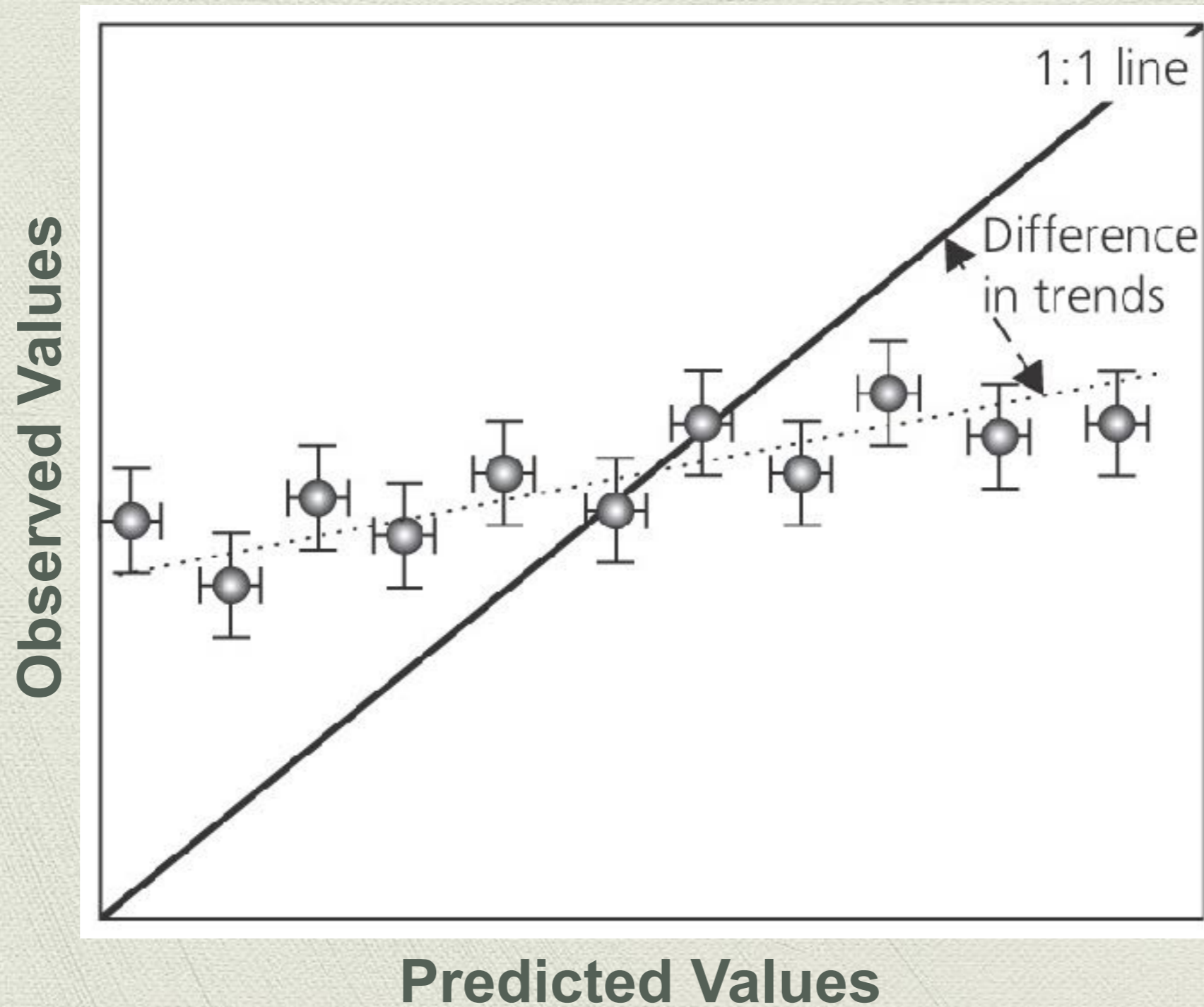
# Accuracy of Prediction

# Identify Outliers

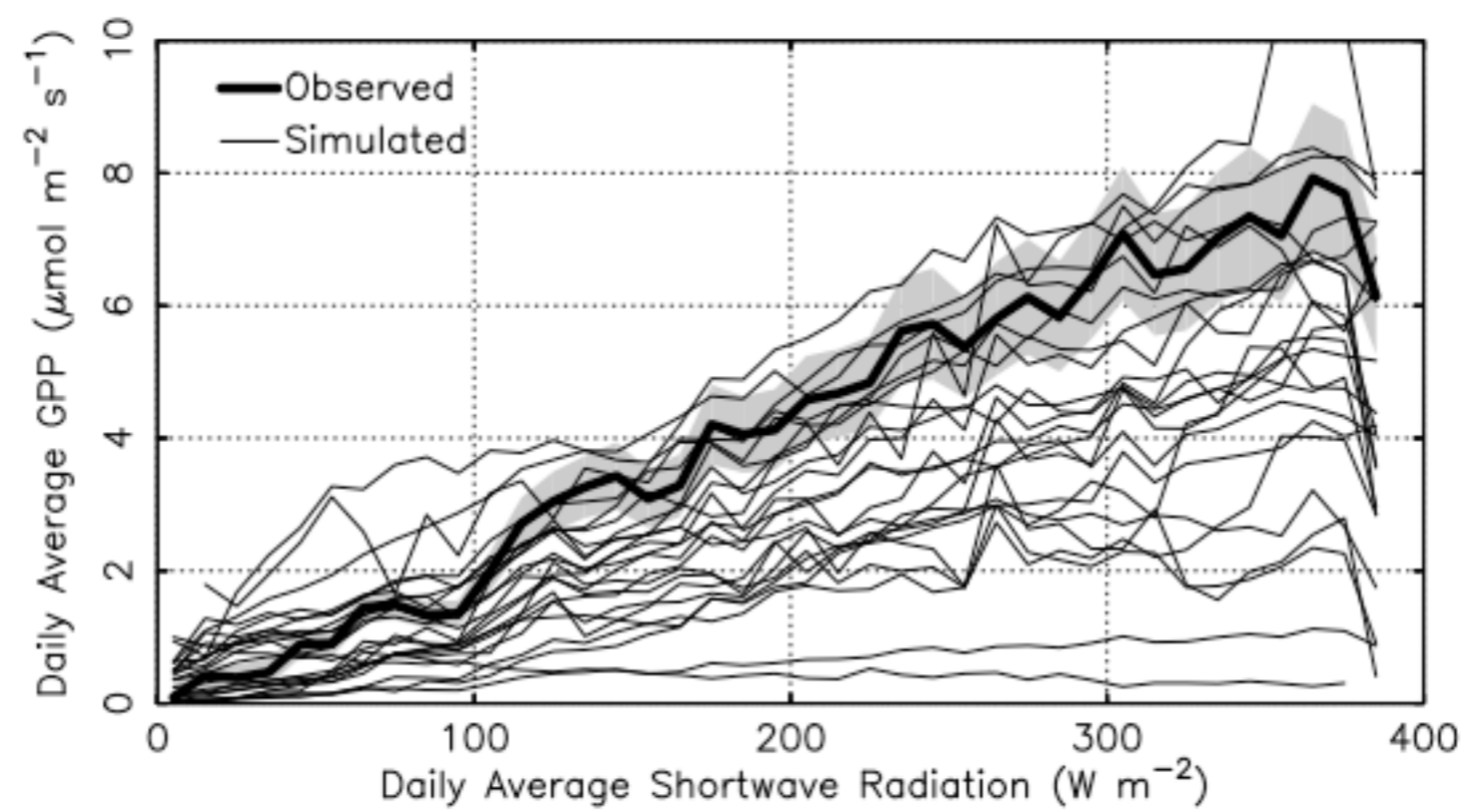# Assess Biases

# Miscalibration

# Dynamics & Drivers

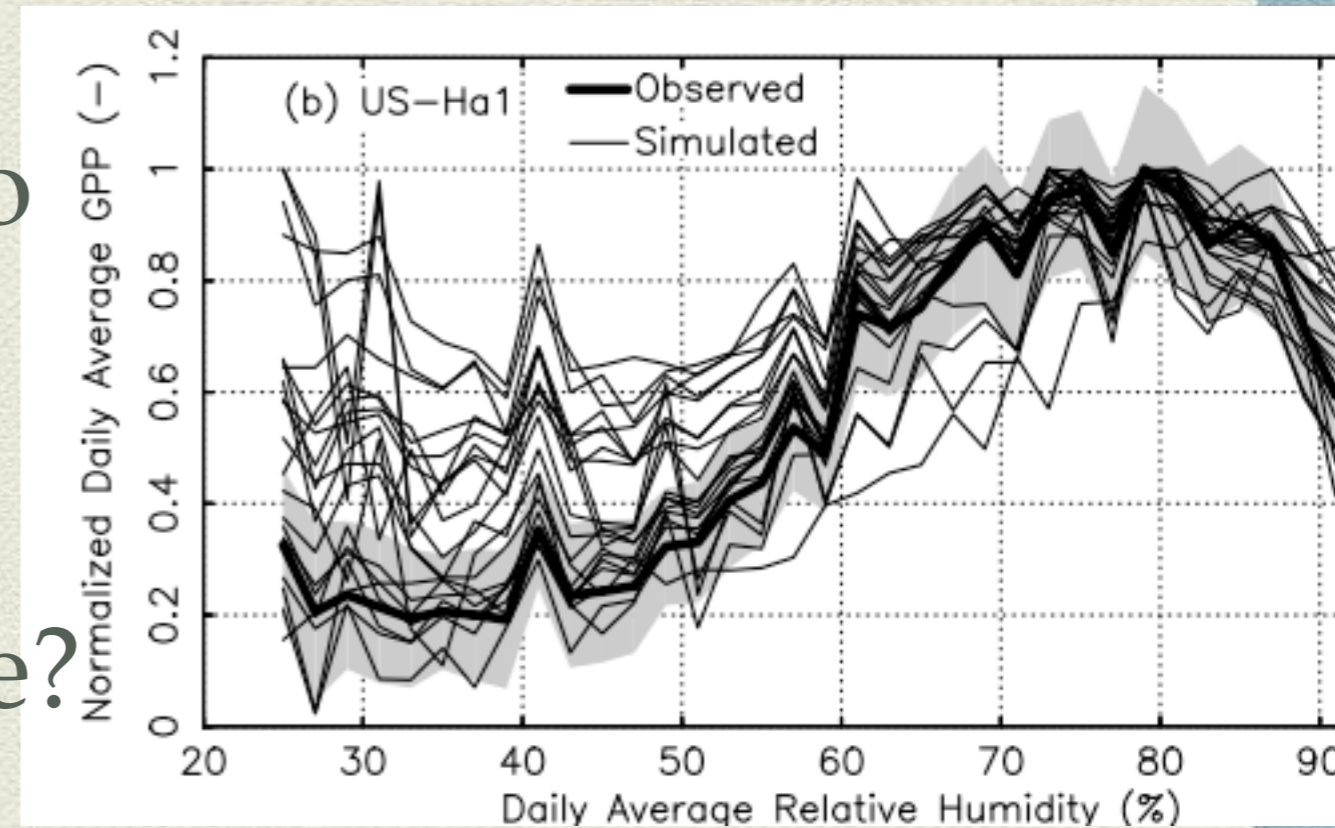Schaefer et al. 2012 JGR-B

# Diagnosing a model is Hypothesis Testing
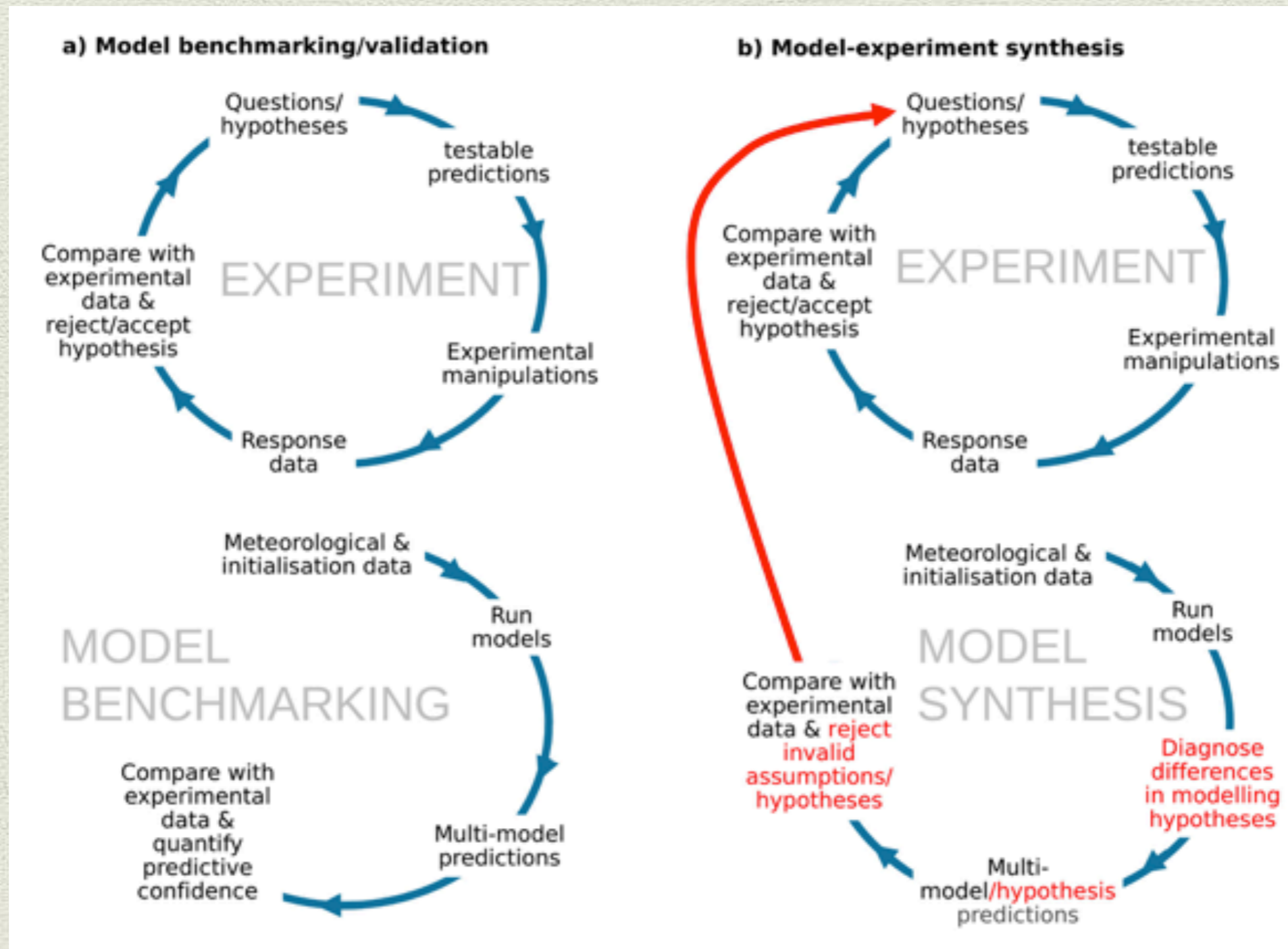
- Why would a model fail at low humidity?
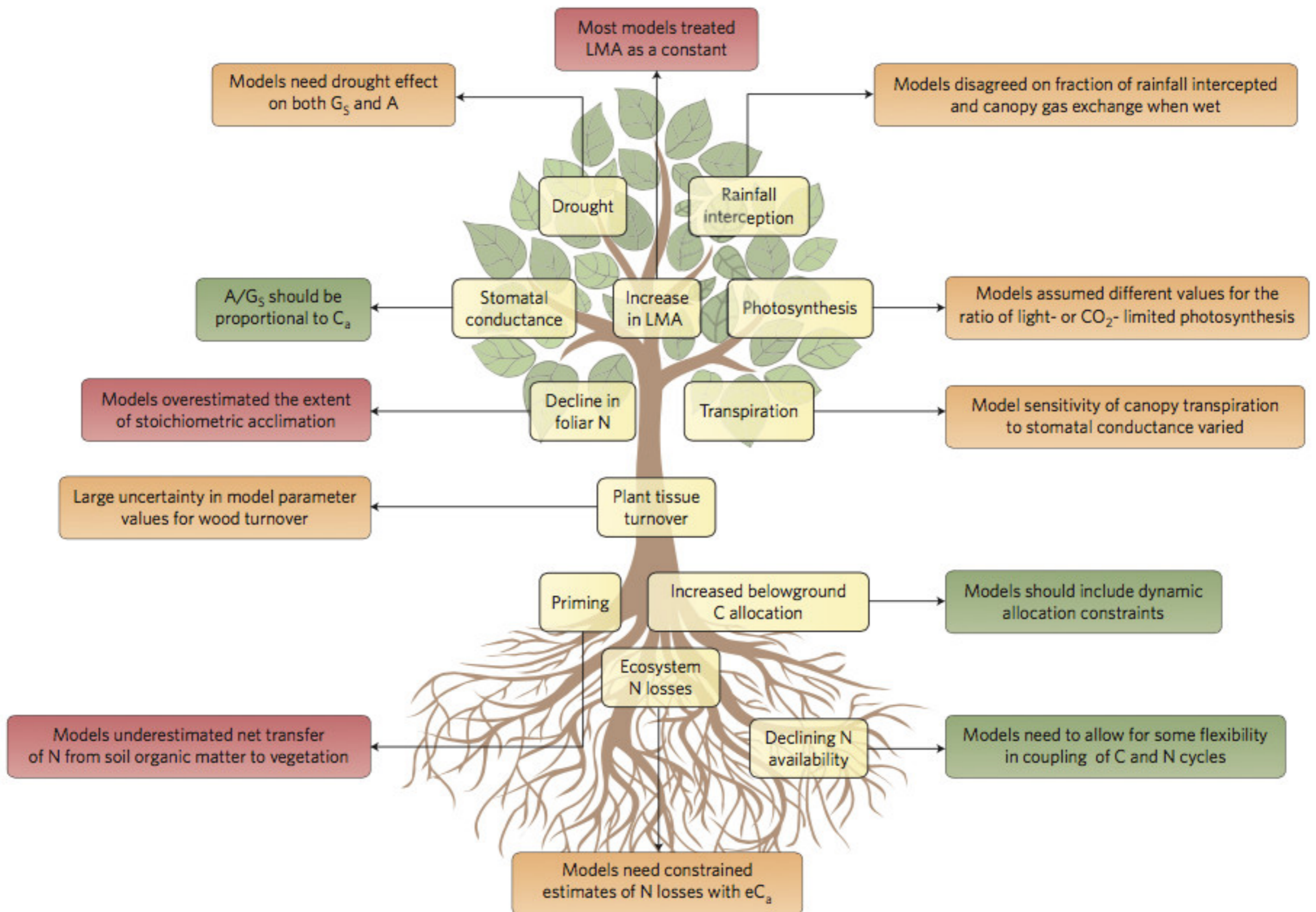
- Stomatal sensitivity too low?

- Too much soil moisture?



- What experiments would I run in the model to test this?

# Focus on key assumptions



Walker et al 2014

Most models treated LMA as a constant

Models need drought effect on both $G_S$ and A

Models disagreed on fraction of rainfall intercepted and canopy gas exchange when wet

Drought

Rainfall interception

$A/G_S$ should be proportional to $C_a$

Stomatal conductance

Increase in LMA

Photosynthesis

Models assumed different values for the ratio of light- or $CO_2$- limited photosynthesis

Models overestimated the extent of stoichiometric acclimation

Decline in foliar N

Transpiration

Model sensitivity of canopy transpiration to stomatal conductance varied

Large uncertainty in model parameter values for wood turnover

Plant tissue turnover

Priming

Increased belowground C allocation

Models should include dynamic allocation constraints

Ecosystem N losses

Models underestimated net transfer of N from soil organic matter to vegetation

Declining N availability

Models need to allow for some flexibility in coupling of C and N cycles

Models need constrained estimates of N losses with $eC_a$
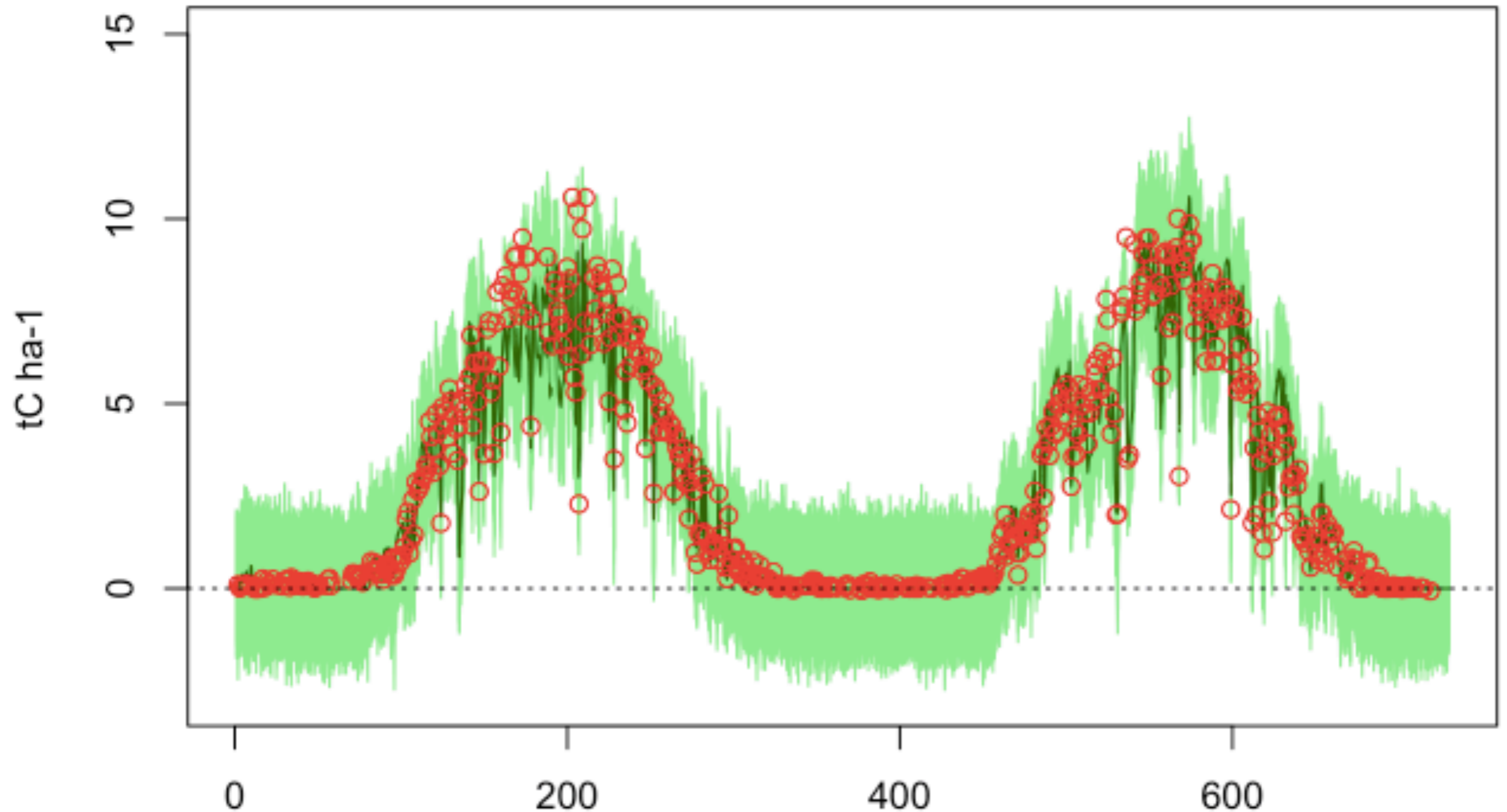
Medlyn et al NCC 2015

"data simulated under a model should look similar to data gathered in the real world."
Conn et al 2018

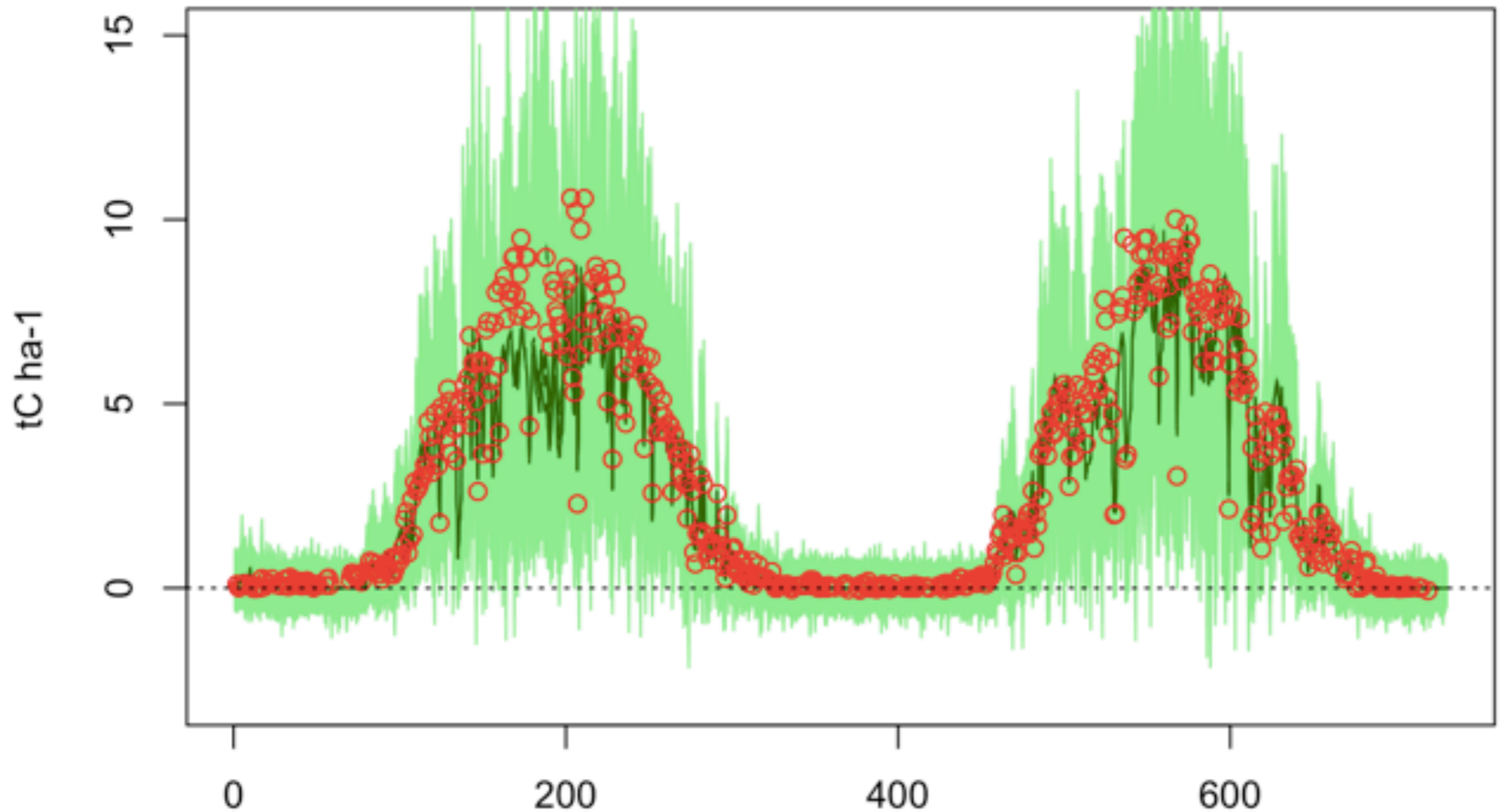# IN THE FITTING, WE ASSUMED IID NORMAL ERRORS



**GPP**

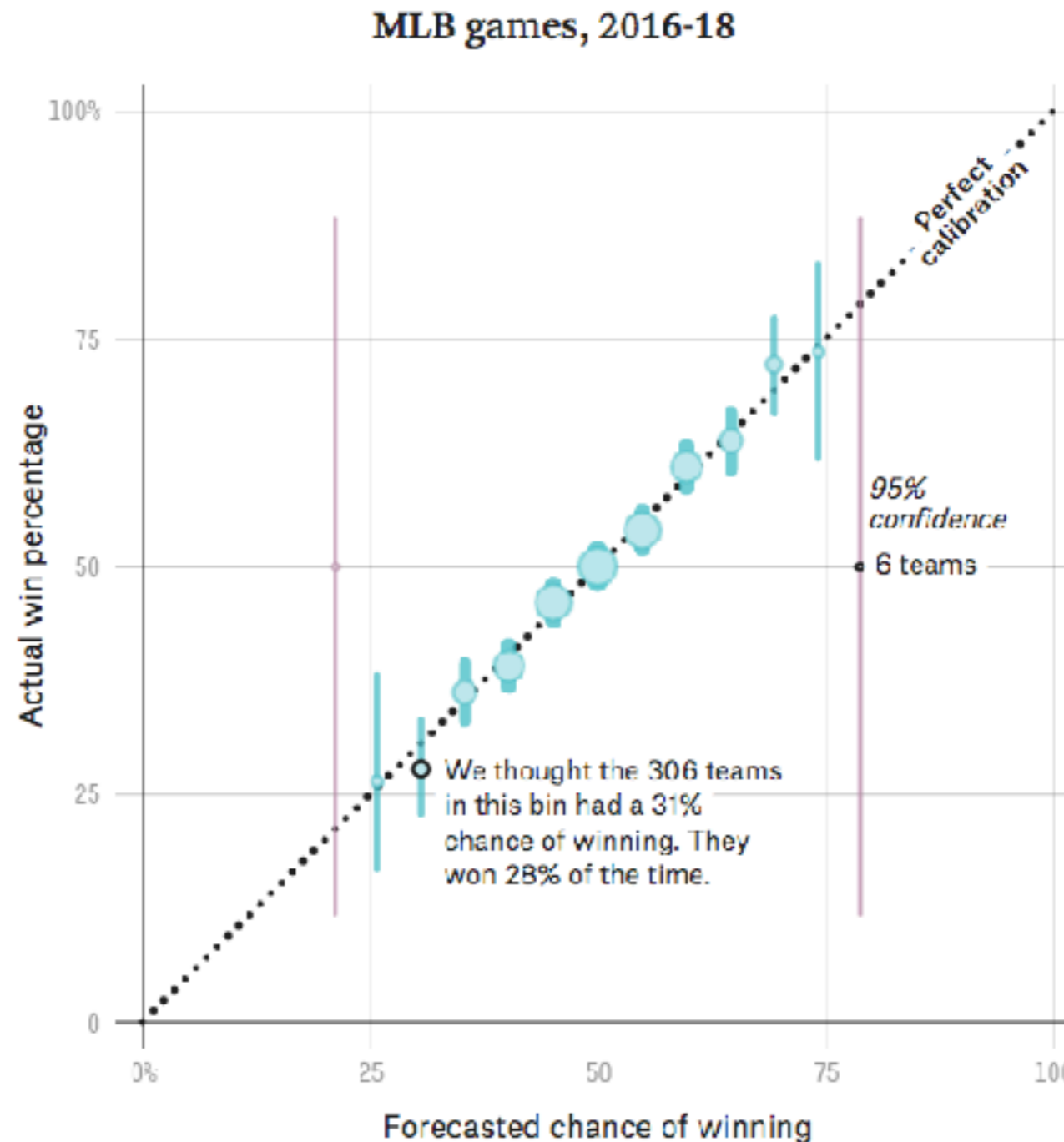Does that seem like an adequate description of the data?

# IN THIS FITTING, WE ASSUMED EXPONENTIAL ERRORS WITH NON-CONSTANT VARIANCE
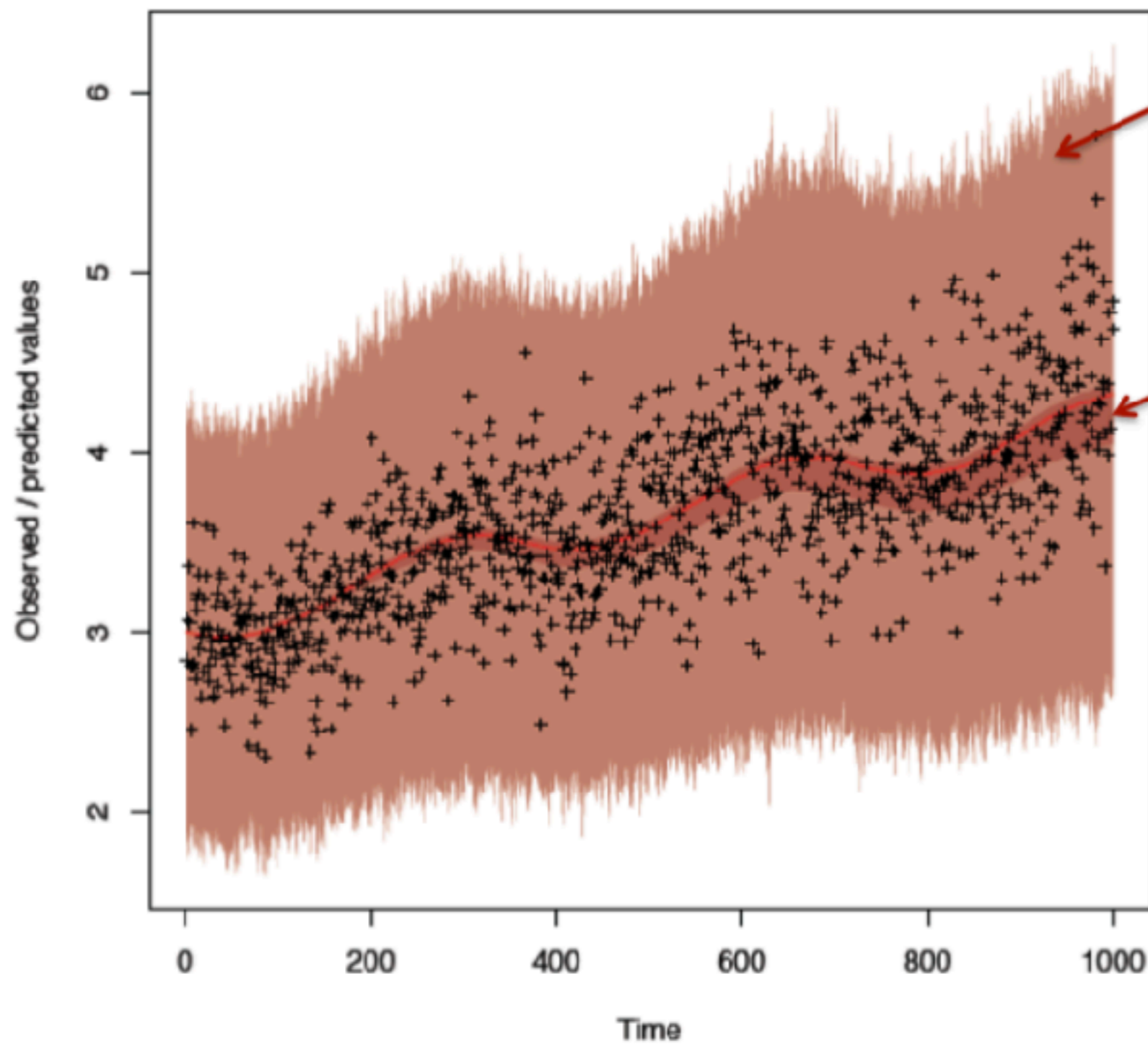


GPP

Does that seem like an adequate description of the data?

# How Good Are FiveThirtyEight Forecasts?



MLB games, 2016-18

We thought the 306 teams in this bin had a 31% chance of winning. They won 28% of the time.

95% confidence

6 teams

Perfect calibration

Actual win percentage

Forecasted chance of winning

https://projects.fivethirtyeight.com/checking-our-work/
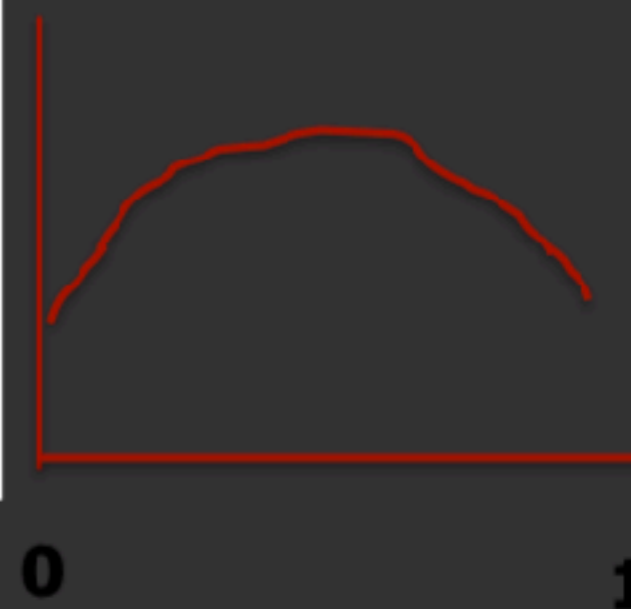
18

# Bayesian p-value / prediction interval

- Posterior predictive distribution is the uncertainty of the „true" value

- **Prediction interval** is the expected variance of the observed values = PPD + error
  - Shows us what distribution we would expect for the data

- Bayesian p-value is when we use PPD + error to calculate the value of the cdf of the observed data
  - Distribution should be flat (uniform)
  - „Bayesian residuals"

PPD + Error

Posterior
Predictive
Distribution

Distribution of ecdf
values for residuals

# Step 3: Quantitative Skill Assessment
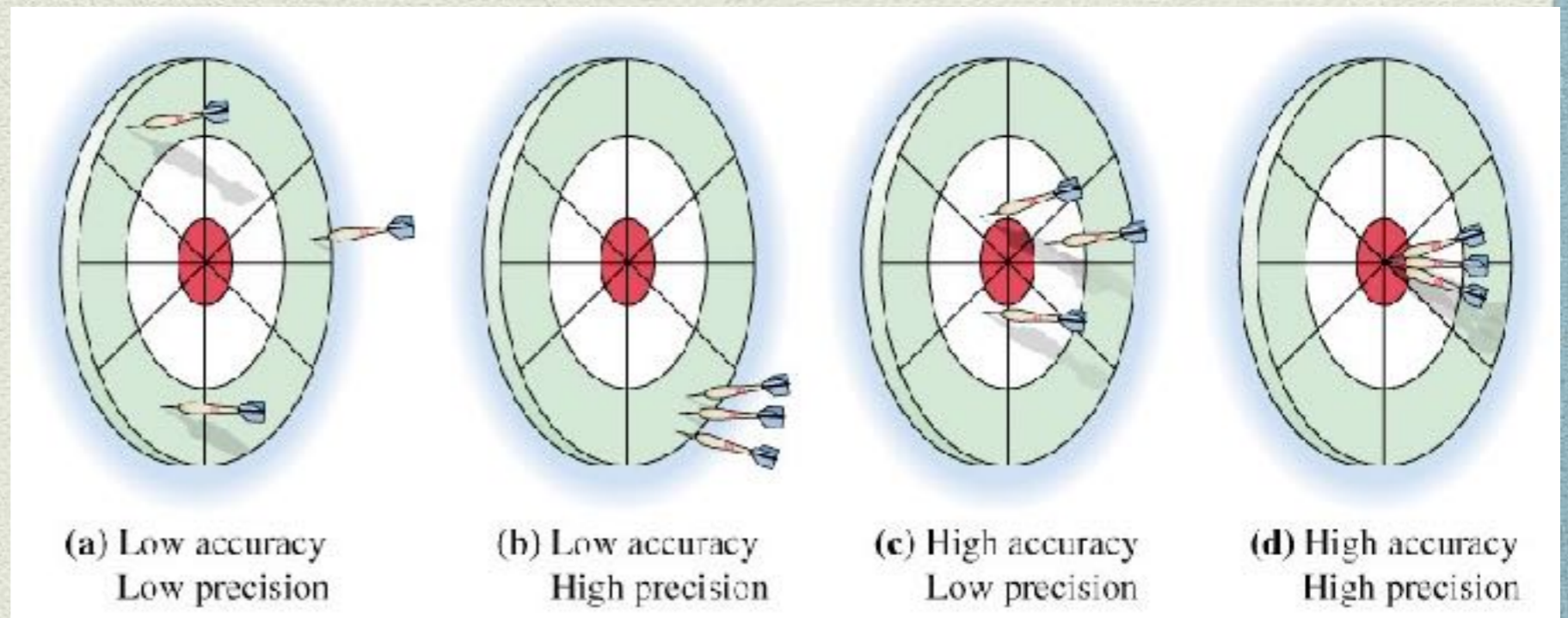
# Error Statistics

- Root Mean Square Error (RMSE)

$$\text{RMSE} = \sqrt{\frac{1}{N}\sum_{i=1}^{N}(y_i - x_i)^2}$$
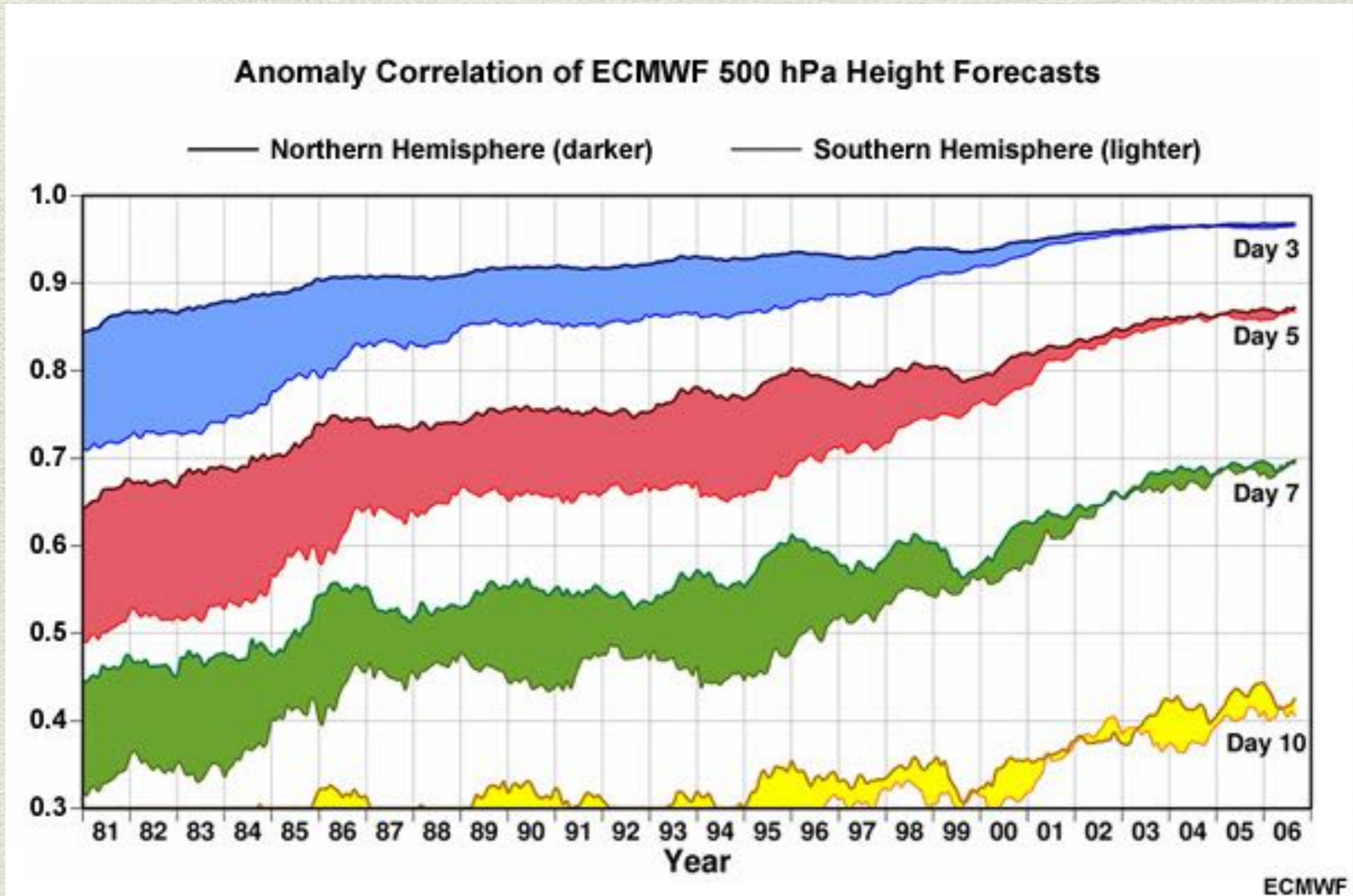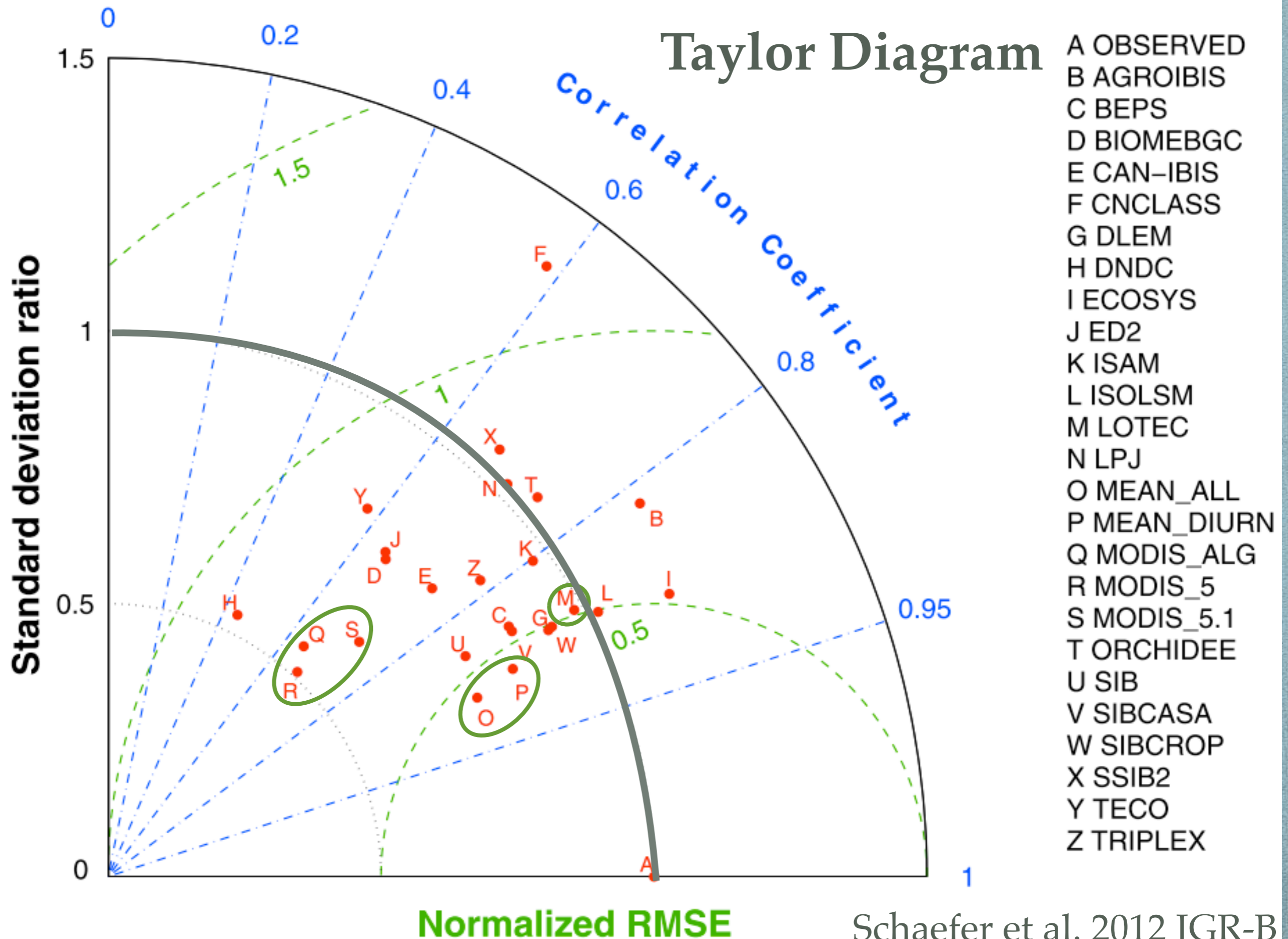
- Bias

- Correlation (r)

- $R^2$

- Regression slope



**(a)** Low accuracy
Low precision

**(b)** Low accuracy
High precision

**(c)** High accuracy
Low precision

**(d)** High accuracy
High precision

**Proper: based on the metric used for calibration**
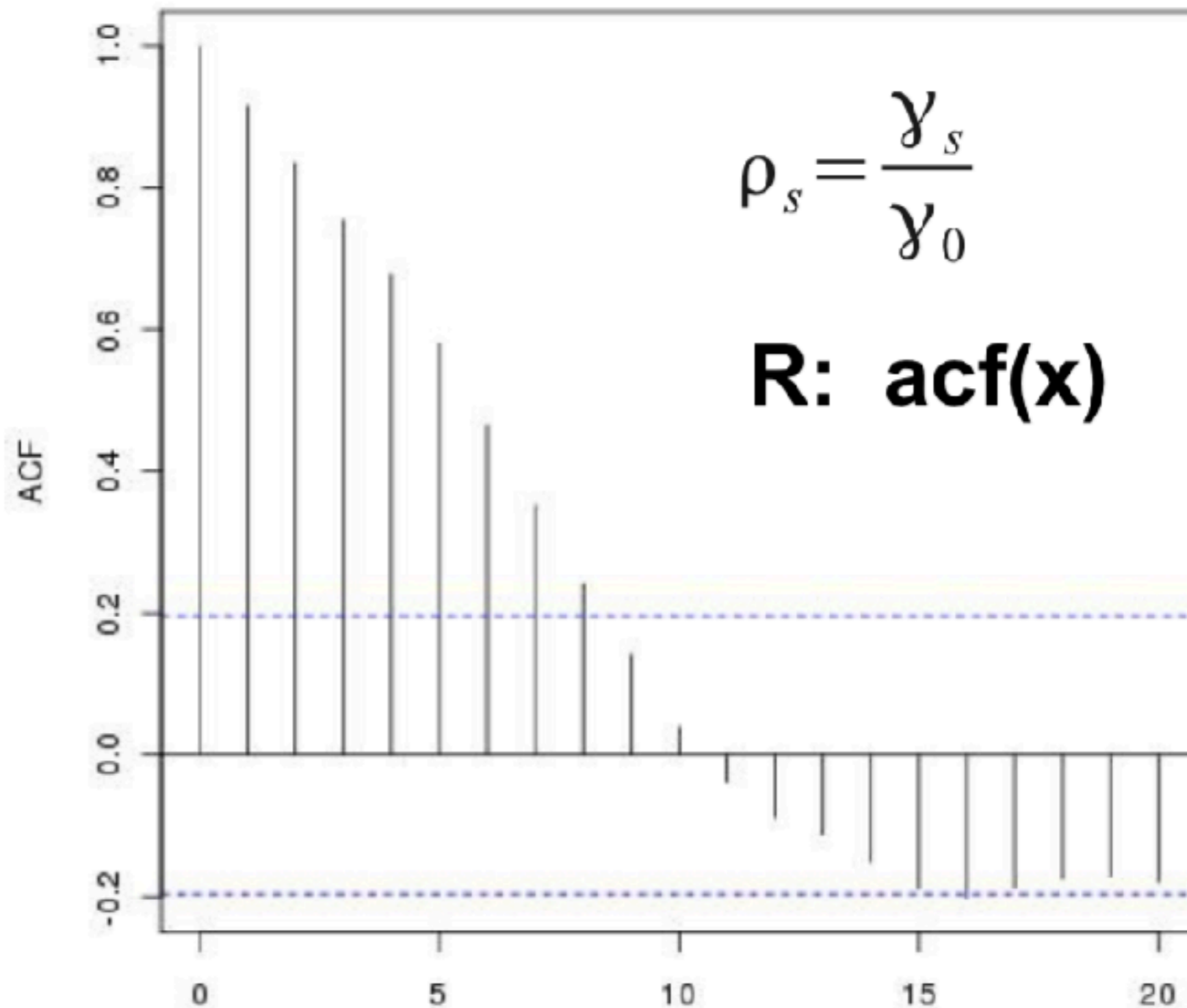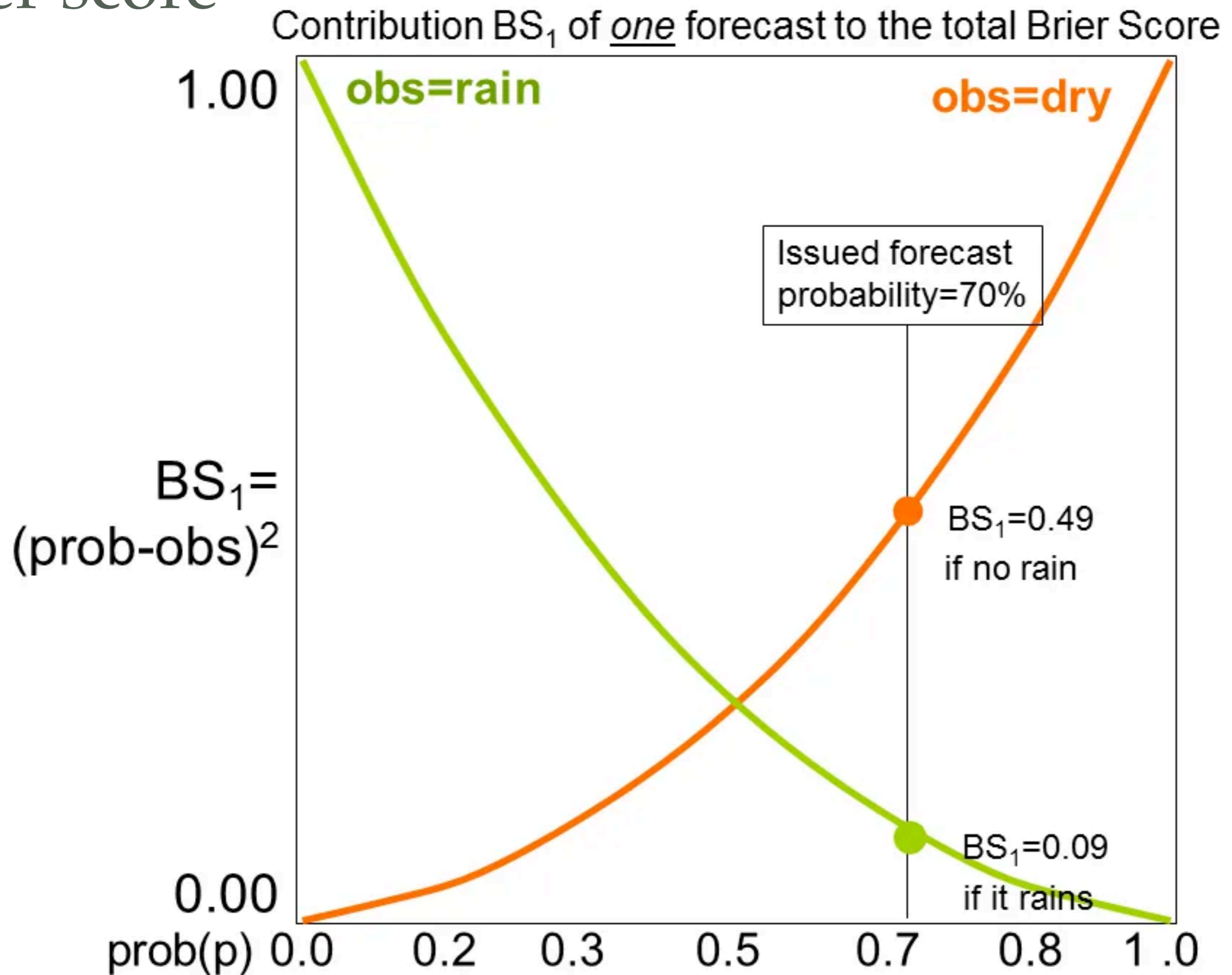**Local: depends on data that could actually be collected**
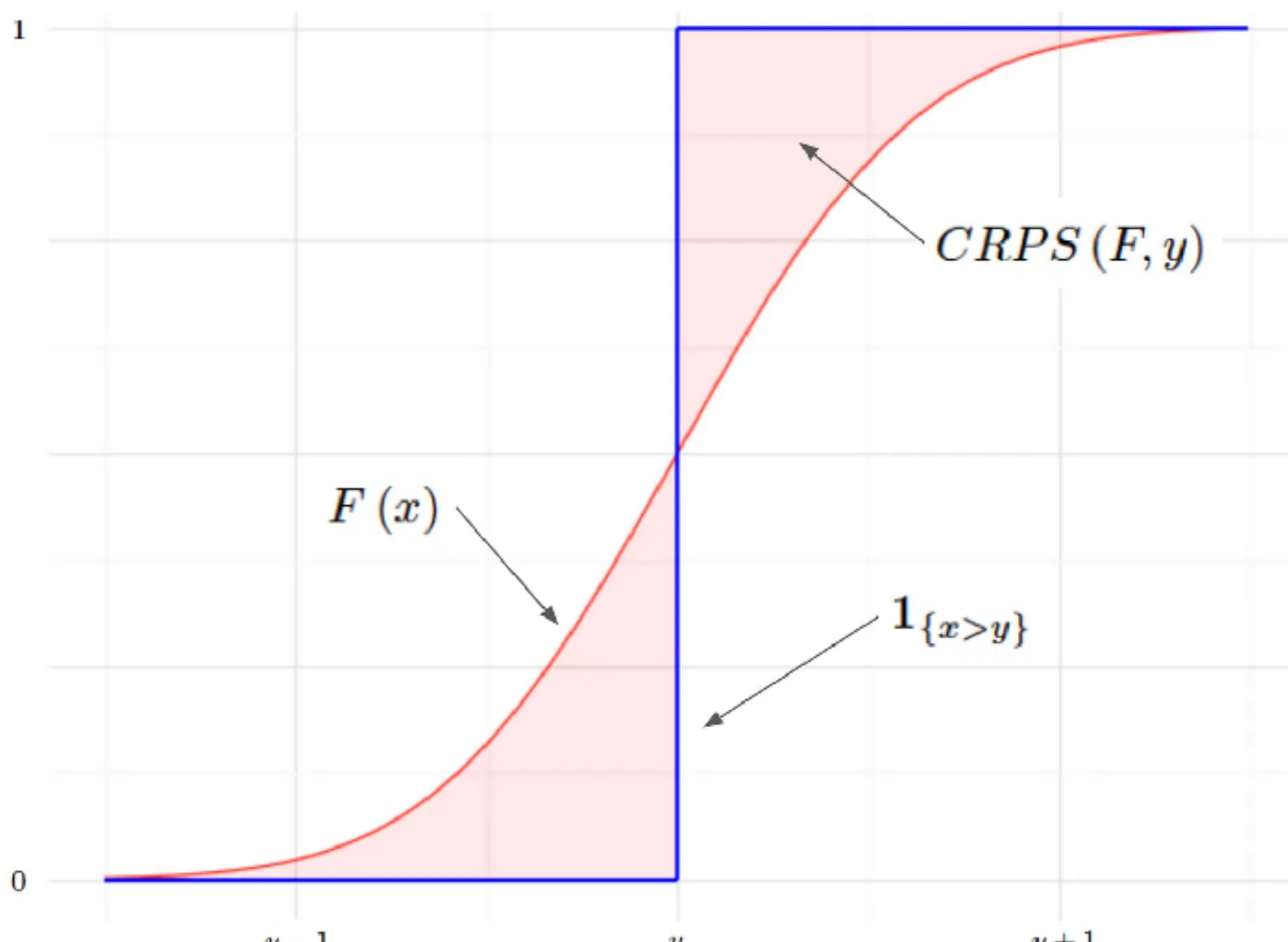
# Correlation



Anomaly Correlation of ECMWF 500 hPa Height Forecasts

Northern Hemisphere (darker)　　Southern Hemisphere (lighter)

Taylor Diagram

A OBSERVED
B AGROIBIS
C BEPS
D BIOMEBGC
E CAN-IBIS
F CNCLASS
G DLEM
H DNDC
I ECOSYS
J ED2
K ISAM
L ISOLSM
M LOTEC
N LPJ
O MEAN_ALL
P MEAN_DIURN
Q MODIS_ALG
R MODIS_5
S MODIS_5.1
T ORCHIDEE
U SIB
V SIBCASA
W SIBCROP
X SSIB2
Y TECO
Z TRIPLEX

Schaefer et al. 2012 JGR-B

# Autocorrelation

## Correlogram



$$\rho_s = \frac{\gamma_s}{\gamma_0}$$

**R: acf(x)**

# Brier score



Contribution $BS_1$ of *one* forecast to the total Brier Score

$BS_1=$
$(prob-obs)^2$

obs=rain

obs=dry

Issued forecast probability=70%

$BS_1=0.49$ if no rain

$BS_1=0.09$ if it rains

prob(p) 0.0    0.2    0.3    0.5    0.7    0.8    1.0

1.00

0.00

# Continuous Ranked Probability Score

$$\mathrm{CRPS}(F, x) = -\int_{-\infty}^{\infty} (F(y) - \mathbb{1}\{y \geq x\})^2 \, \mathrm{d}y$$

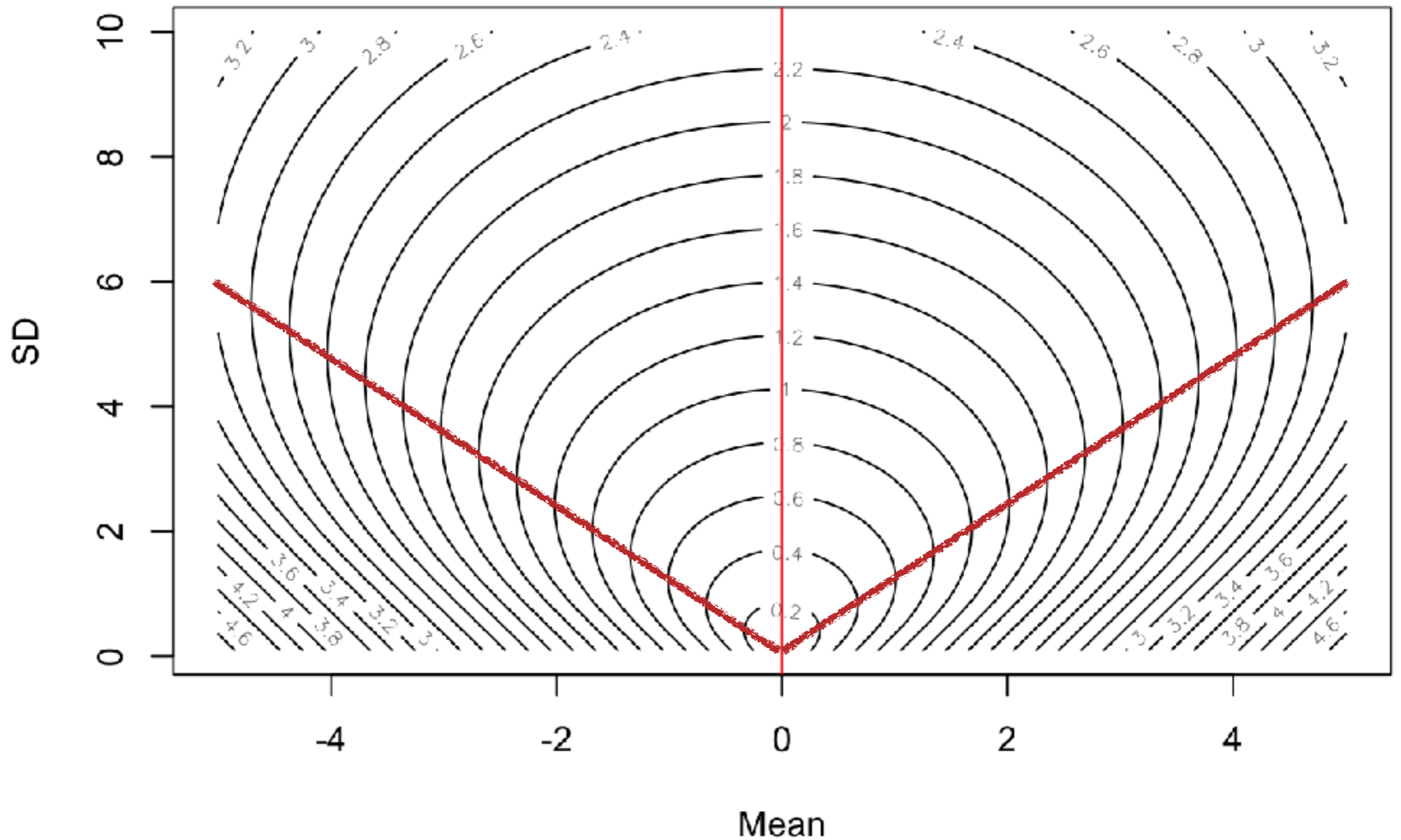# Continuous Ranked Probability Score

$$\text{CRPS}(F, x) = -\int_{-\infty}^{\infty} (F(y) - \mathbb{1}\{y \geq x\})^2 \, \mathrm{d}y$$

**Ensemble member**

**Data**

$$\text{CRPS}(\hat{F}_m, y) = \frac{1}{m}\sum_{i=1}^{m} |X_i - y| - \frac{1}{2m^2}\sum_{i=1}^{m}\sum_{j=1}^{m} |X_i - X_j|$$

**Mean Absolute Error**

**Penalty for ensemble spread**

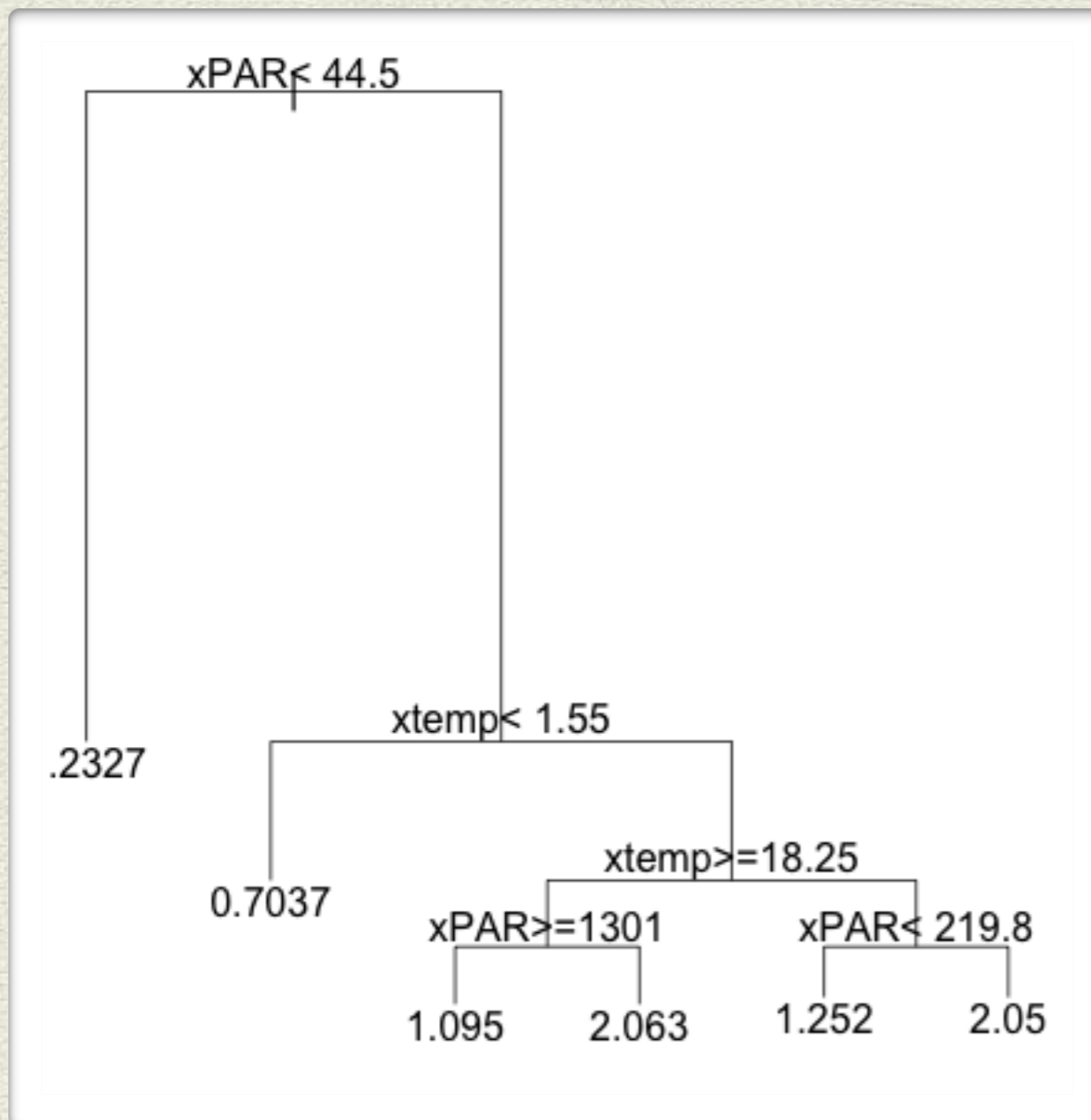https://github.com/eco4cast/neon4cast-scoring/blob/main/
CRPS_example_JRT.Rmd

# Data mining the residuals

- Wide variety of Data Mining algorithms in use

- Large debate about use in process modeling and forecasting

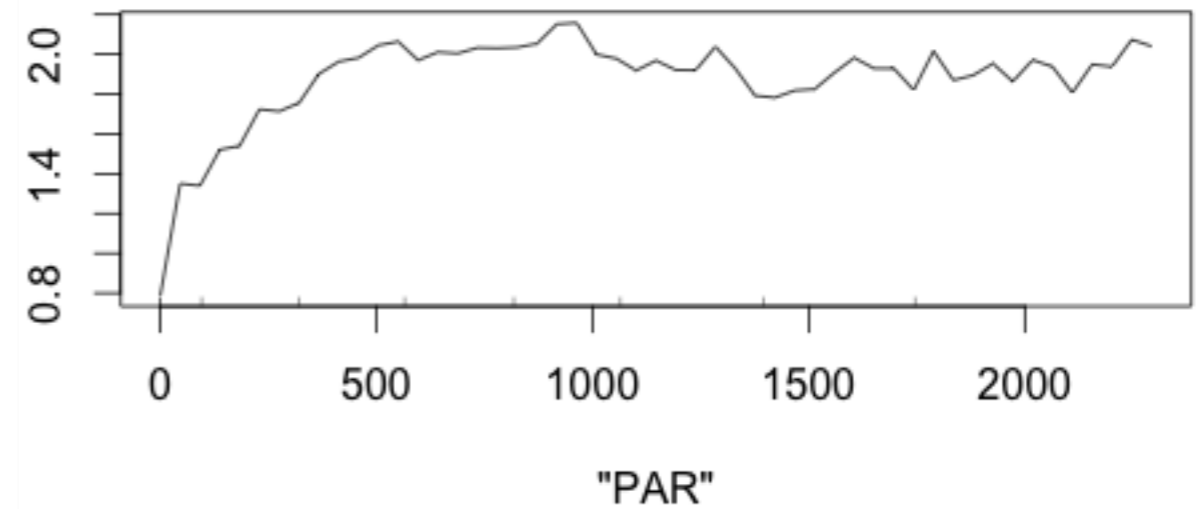- Potentially useful for generating hypothesis about when/where model fails

- CART

- GAM

- Random Forests

- Boosted regression trees / XGBoost

- Artificial Neural Network

- Deep Learning

Hybrid Models (Process + NN)
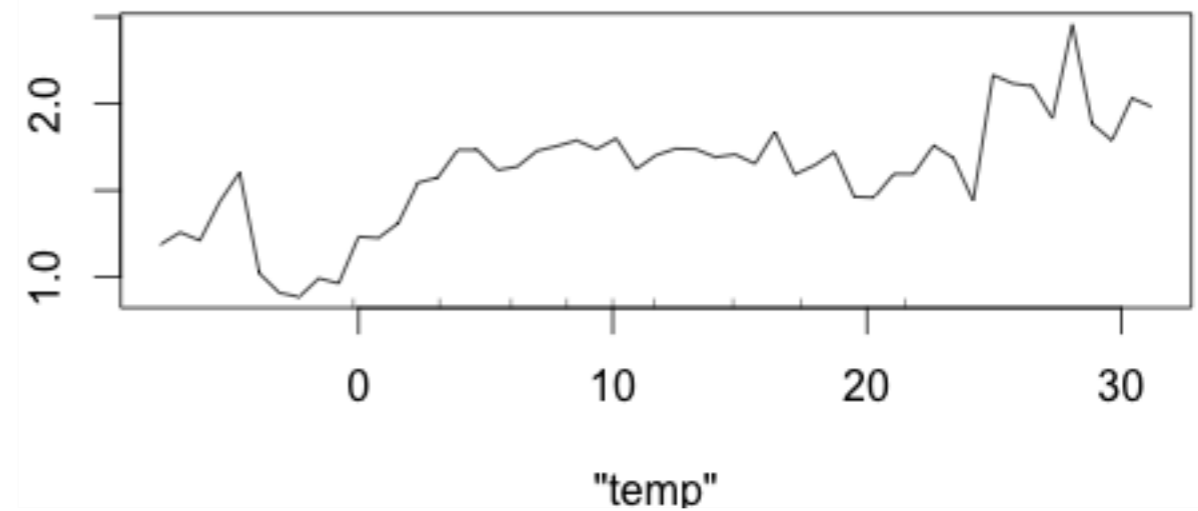
**A. Naïve Neural Network**

Input Layer — Hidden Layers — Output Layer

Environmental Covariates $X \rightarrow$ Environmental system state $Y$

**B. Bias Correction**

Input Layer — Hidden Layers — Output Layer

$X \rightarrow$ Process-based model $\hat{y}_{\text{PHY}} \rightarrow Y$

Parameters $\tau$

**C. Parallel Physics**

env covariates explain bias

Input Layer — Hidden Layers — Output Layer

$X \rightarrow Y$

Process-based model $\hat{y}_{\text{PHY}}$

Parameters $\tau$

**D. Physics Regularisation**

Input Layer — Hidden Layers — Output Layer

$X \rightarrow Y$

Process-based model $\hat{y}_{\text{PHY}}$

weight given to process model

Parameters $\tau$

**E. Domain Adaptation**

NN pretrained on model outputs

Input Layer — Hidden Layers pretrained — Hidden Layers retrained — Output Layer

then retrained against data

$X \rightarrow Y$

$X_{\text{SIM}}$ — Process-based model $\hat{y}_{\text{PHY}}$

Parameters $\tau$

**F. Physics Embedding**

Input Layer — Parameter Net $\tau$ — Residual Net — Output Layer

$X \rightarrow$ Process-based model $\hat{y}_{\text{PHY}} \rightarrow Y$

NN predicts parameter variability

Bias Correction

Wesselkamp et al 2024 Ecology Letters
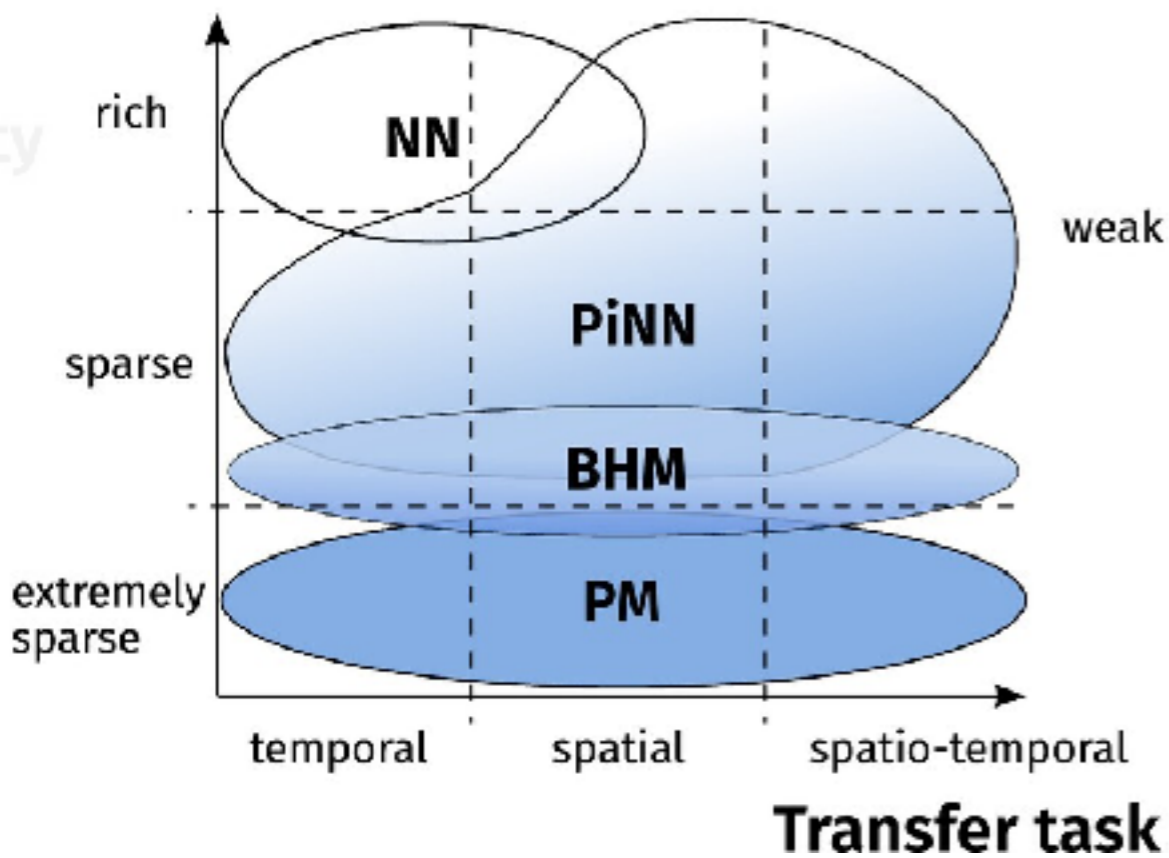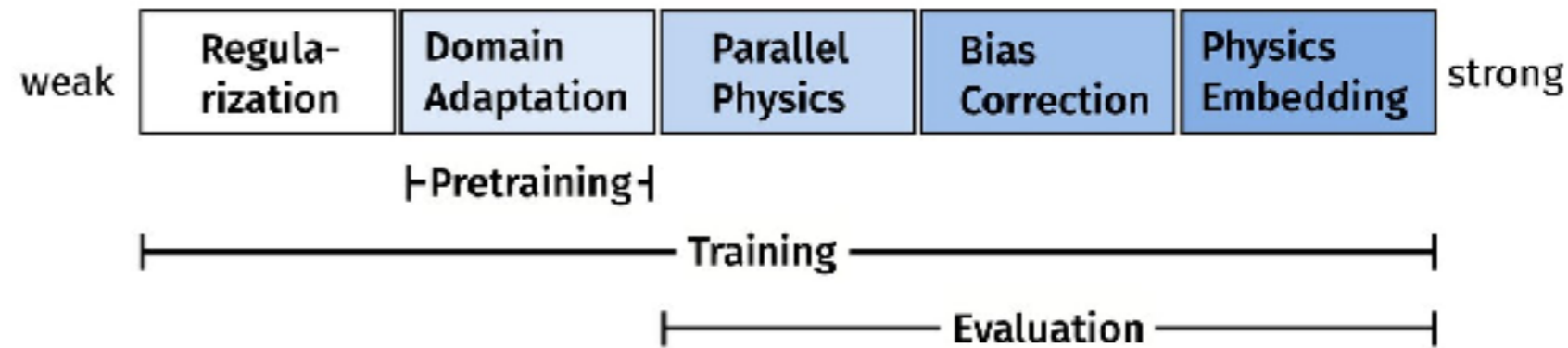
Expected performance sweet spots

Data Availability

Theory constraint of PiNNs

Wesselkamp
et al 2024
Ecology Letters