

Improved return level estimates of cyclone-induced extreme waves:

combining extreme value distribution and probabilistic machine learning predictions

Jeremy Rohmer*, Andrea Filippini*, Rodrigo Pedreros*,
*BRGM, French Geological survey j.rohmer@brgm.fr

Karina

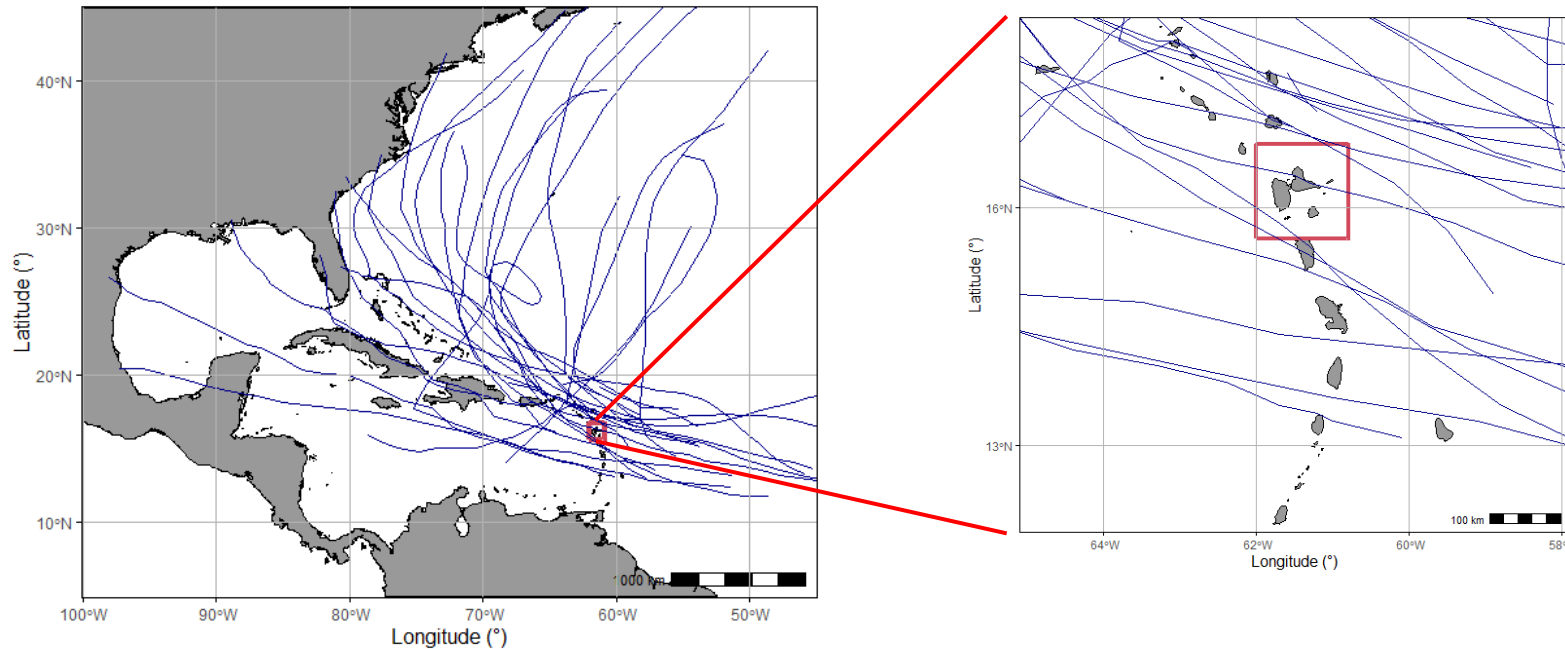
Teddy

Vicky



September 14, 2020, NOAA's GOES-East satellite spied six active tropical systems spanning the Atlantic and Pacific Ocean - <https://www.nesdis.noaa.gov/content/six-tropical-systems-swirl-around-two-oceans>

Motivating test case: cyclone-induced extreme waves at Guadeloupe archipelago (French West Indies)



Objective: estimate extreme wave heights (H_s) induced by tropical cyclones TC along the Guadeloupe coast

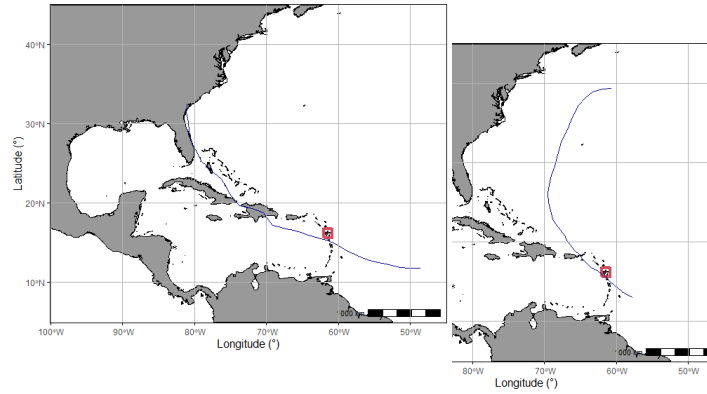
Difficulty: historical data are relatively scarce (28 TCs)

 Cyclone Track (HURDAT database 1979-2019) at a distance <400km from the Island

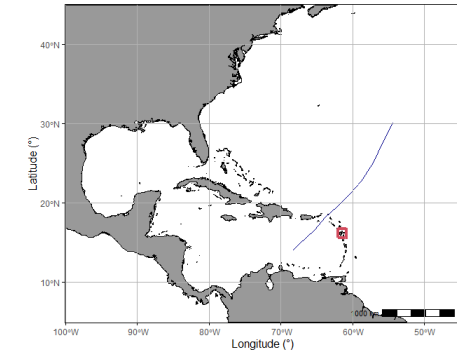
Synthetic TC approach

1

Randomly generate a large number of synthetic TCs using stochastic generator [1] **~700 TCs (equivalent of 1,000 years)**



...



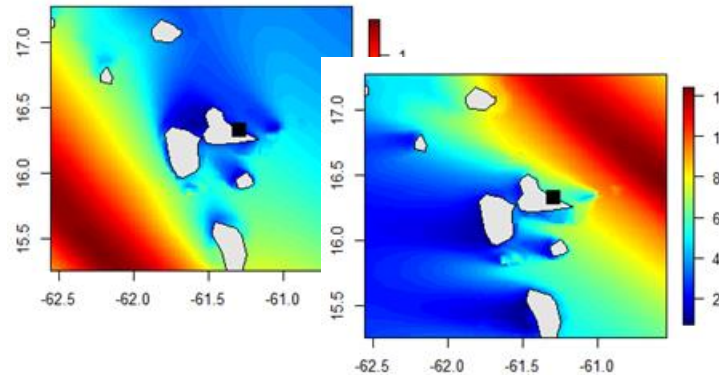
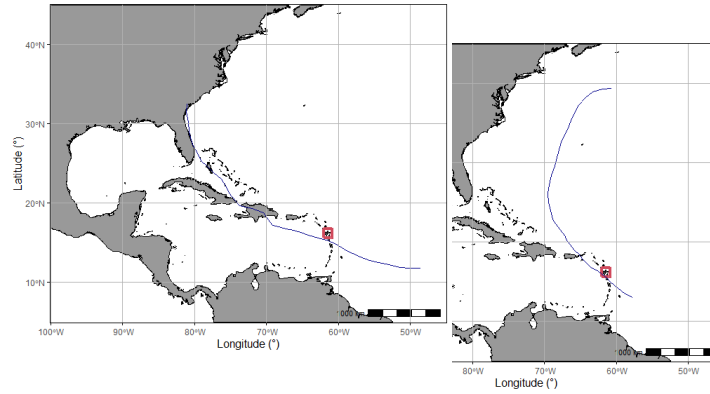
Synthetic TC approach

1

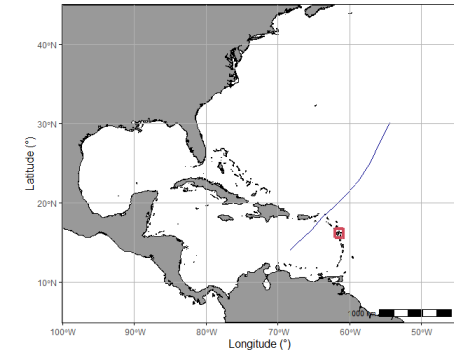
Randomly generate a large number of synthetic TCs using stochastic generator [1] **~700 TCs (equivalent of 1,000 years)**

2

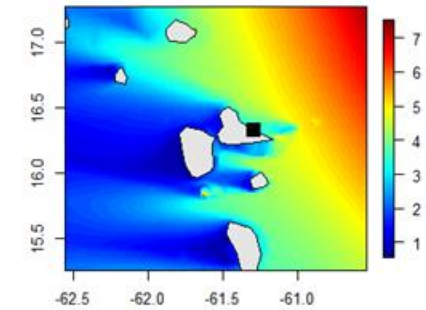
Compute the TC-induced wave heights using numerical hydrodynamic simulators (here WW3 version 4.18 [2])



...



...



Synthetic TC approach

1

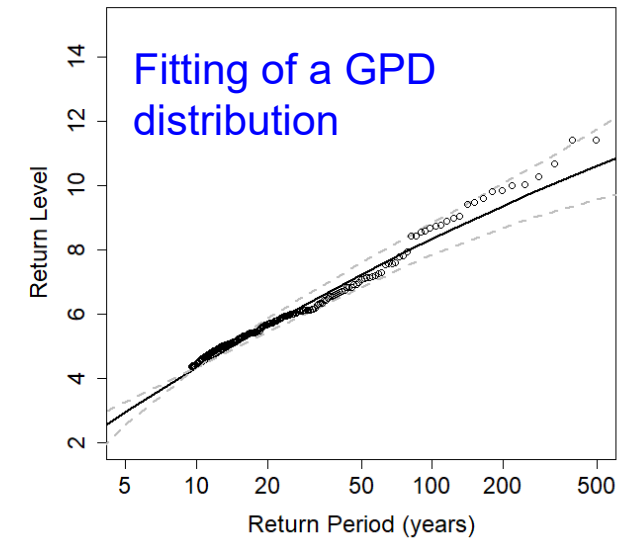
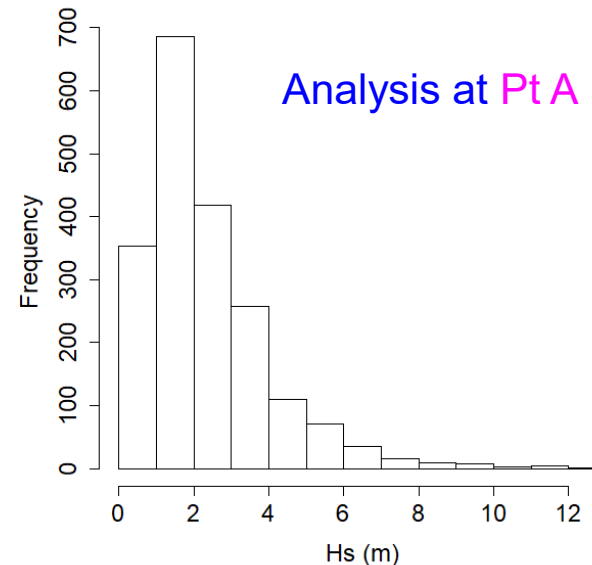
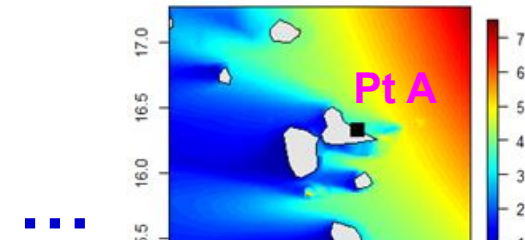
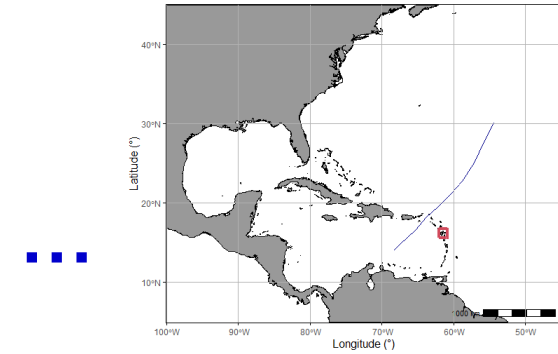
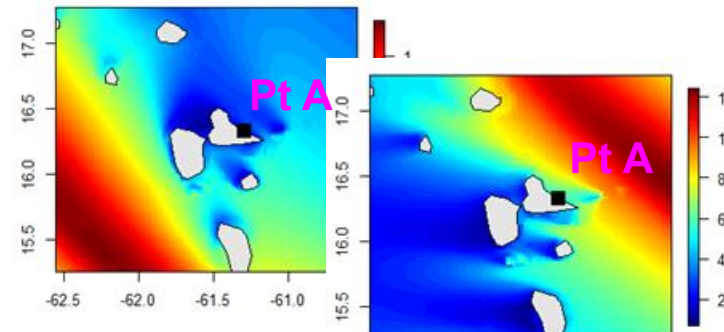
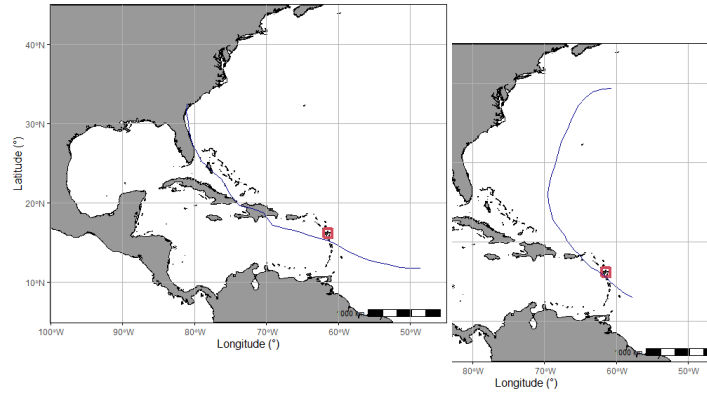
Randomly generate a large number of synthetic TCs using stochastic generator [1] **~700 TCs (equivalent of 1,000 years)**

2

Compute the TC-induced wave heights using numerical hydrodynamic simulators (here WW3 version 4.18 [2])

3

Extract Hs maximum values around the coasts of Guadeloupe and perform **extreme value analysis** [3] to evaluate the return levels of interest (e.g. 100y Hs RL)

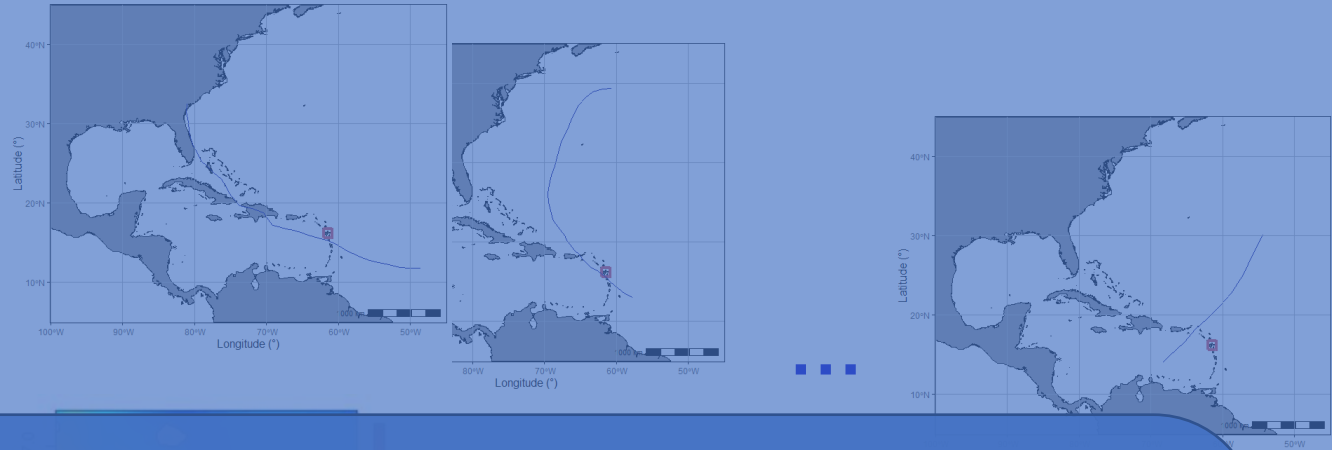


[1]: Bloemendaal et al. (2020); [2]: Tolman (2014); [3]: Coles et al., 2001

Synthetic TC approach

1

Randomly generate a large number of synthetic TCs using stochastic generator [1] ~700 TCs (equivalent of 1 000 years)



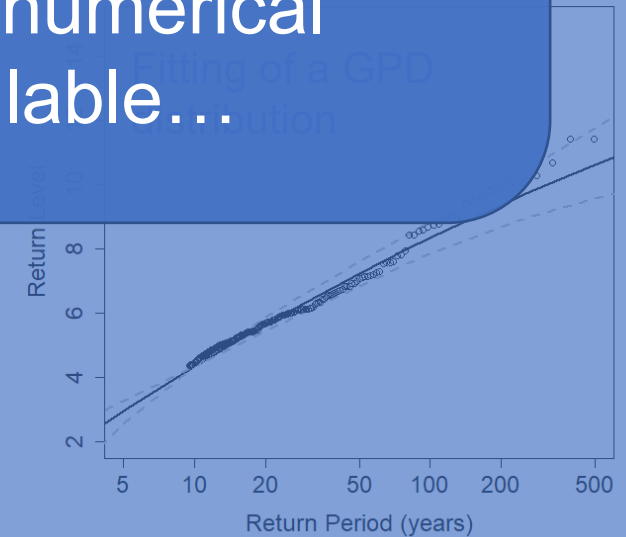
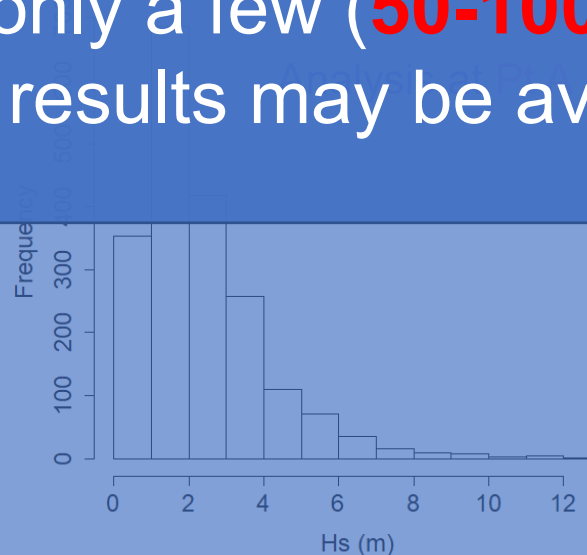
2

Compute the heights using hydrodynamic WW3 version 4.18 [2]

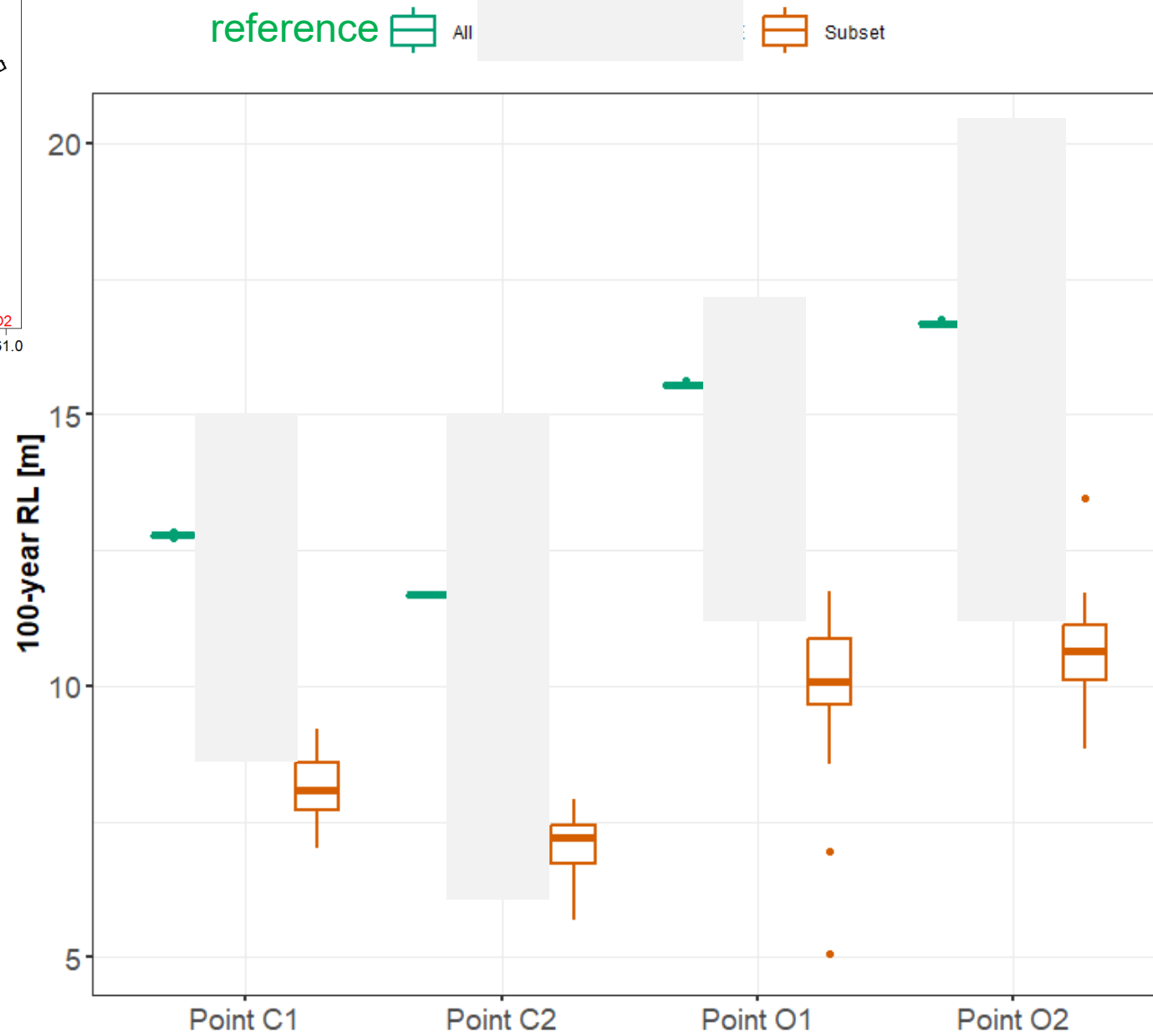
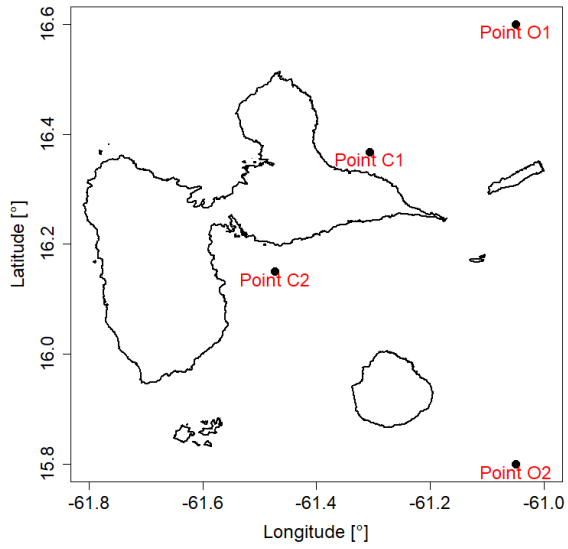
The computational time cost of each numerical simulation is large (several hours, even days)

3

Extract around the coasts of Guadeloupe and perform extreme value analysis [3] to evaluate the return levels of interest (e.g. 100y Hs RL)



Subset of p=10% (70 cyclones)



Use a machine learning (ML) models?

1

Train a ML model with a **subset of cyclones** to model the relationship between TCs characteristics (radius, atmospheric pressure, distance to cyclone eye) and the waves >> **here random forest (RF) [1]**

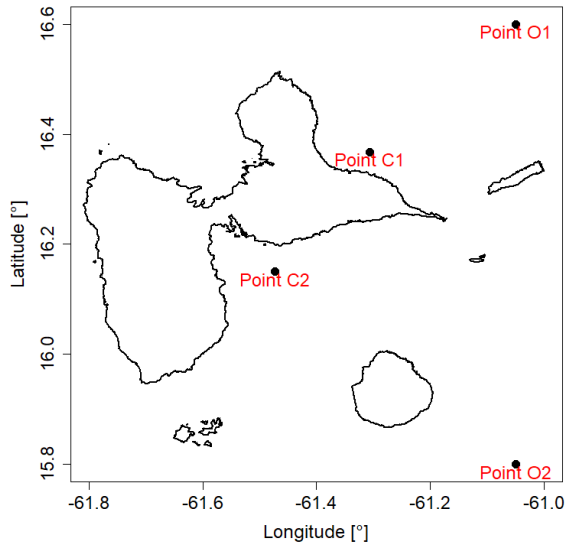
2


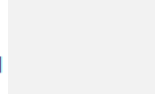


Use the RF model to **predict waves** for a larger number of cyclones >> here for the 700 synthetic cyclones

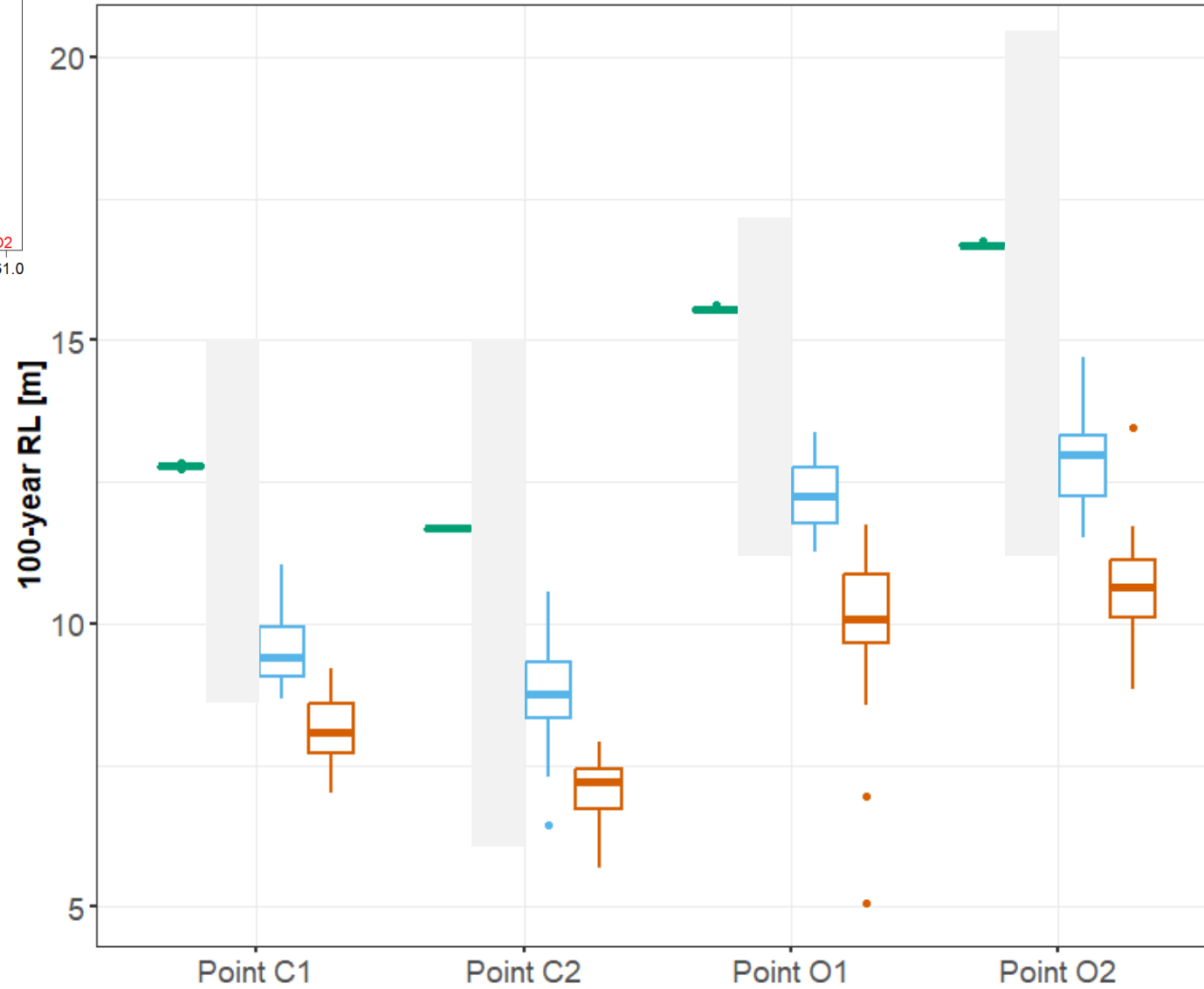
3

Estimate the 100y RL on **the augmented dataset**

'Crude' ML approach 'RFwoE'

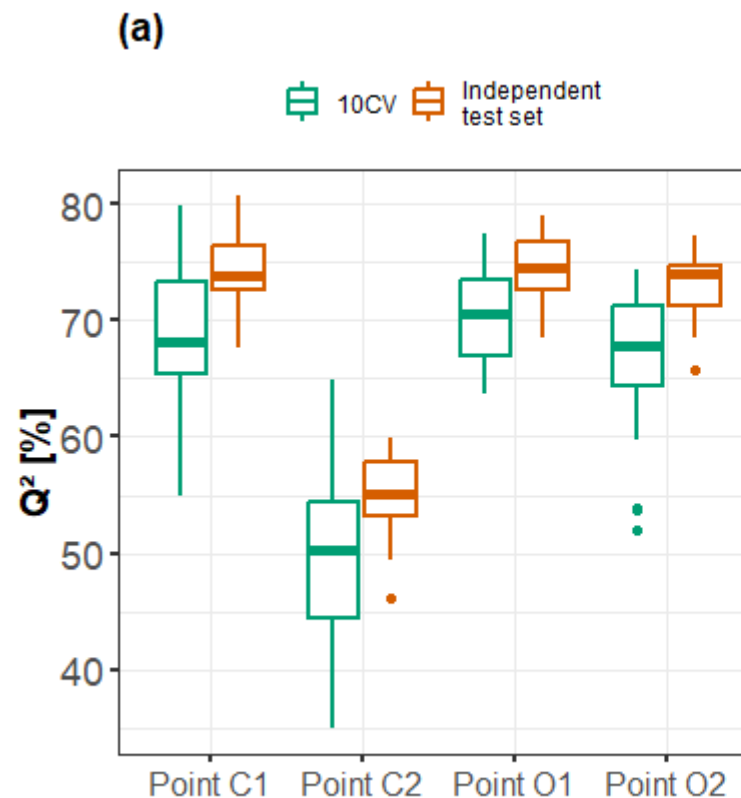


reference  All  RFwoE  Subset 

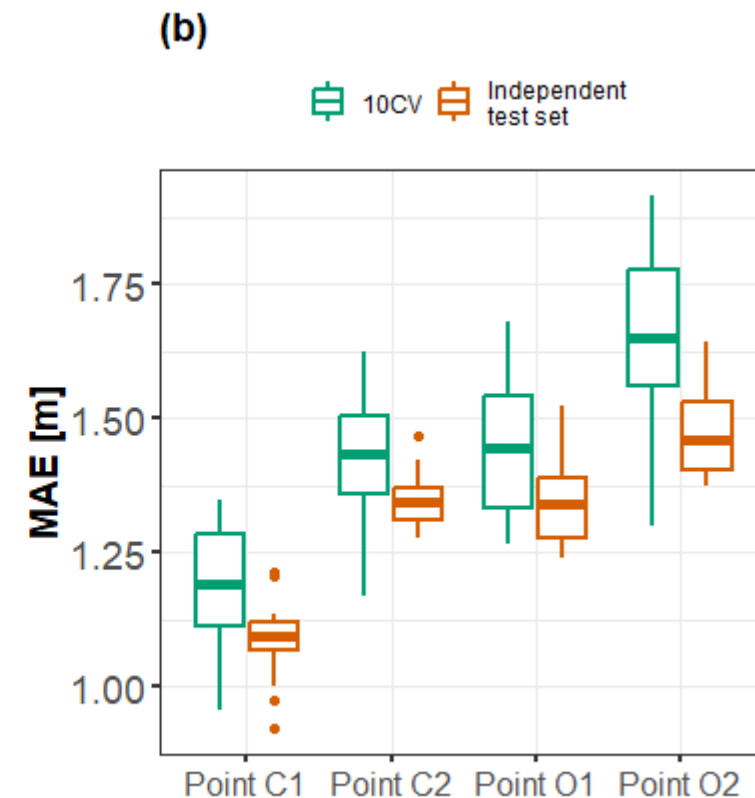


ML does not provide a 'perfect' prediction ($Q^2 \sim 80\%$, MAE $\sim 1\text{m}$)

IDEA: reflect the prediction error in the RL estimates



The higher, the better



The lower, the better

In practice,

- Reflect the **prediction error via a probabilistic ML** >> quantile RF [1]
 - Account for the prediction error along the inference of the parameters θ of the extreme value probability distribution
- >> use of the flexible **Approximate Bayesian Computation** scheme [1] with **Wasserstein W distance** [2]

Flexible Bayesian inference scheme

1. Avoids the computation of the likelihood
2. Easily integrated any ML probabilistic distribution via sampling

Robust metric for measuring the distance between probability distributions

In practice,

- Reflect the **prediction error via a probabilistic ML** >> quantile RF [1]
 - Account for the prediction error along the inference of the parameters θ of the extreme value probability distribution
- >> use of the flexible **Approximate Bayesian Computation** scheme [1] with **Wasserstein W distance** [2]
- 1 Draw θ from the prior probability distributions: $\theta \sim p(\theta)$;

In practice,

- Reflect the **prediction error via a probabilistic ML** >> quantile RF [1]
 - Account for the prediction error along the inference of the parameters θ of the extreme value probability distribution
- >> use of the flexible **Approximate Bayesian Computation** scheme [1] with **Wasserstein W distance** [2]

1 Draw θ from the prior probability distributions: $\theta \sim p(\theta)$;

2 Draw a n_p -dimensional vector $\tilde{U} = (\tilde{u}_1, \dots, \tilde{u}_{n_p})$ of uniformly distributed values between 0 and 1;

3 Simulate an n_p -dimensional vector of random realisations $\tilde{Y}^{\text{meta}} = (q^{\tilde{u}_1}(y|\mathbf{x}^*), \dots, q^{\tilde{u}_{n_p}}(y|\mathbf{x}^*))$ using the procedure described in Set. 2.2 (based on the qRF model fitted using n training data);

Propagate the prediction error

In practice,

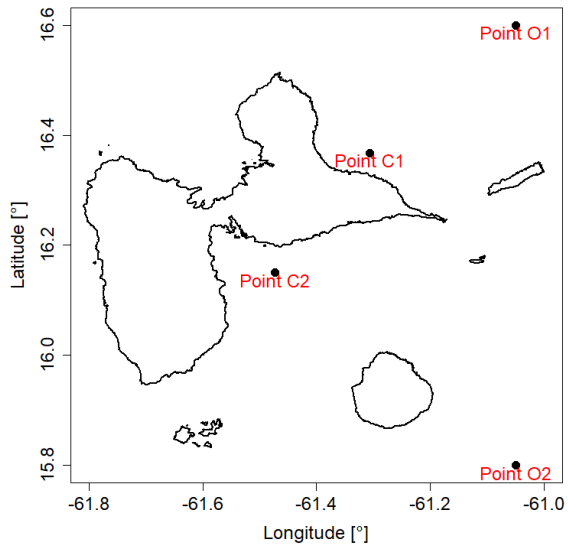
- Reflect the **prediction error via a probabilistic ML** >> quantile RF [1]
 - Account for the prediction error along the inference of the parameters θ of the extreme value probability distribution
- >> use of the flexible **Approximate Bayesian Computation** scheme [1] with **Wasserstein W distance** [2]

- 1 Draw θ from the prior probability distributions: $\theta \sim p(\theta)$;
- 2 Draw a n_p -dimensional vector $\tilde{U} = (\tilde{u}_1, \dots, \tilde{u}_{n_p})$ of uniformly distributed values between 0 and 1;
- 3 Simulate an n_p -dimensional vector of random realisations $\tilde{Y}^{\text{meta}} = (q^{\tilde{u}_1}(y|\mathbf{x}^*), \dots, q^{\tilde{u}_{n_p}}(y|\mathbf{x}^*))$ using the procedure described in Set. 2.2 (based on the qRF model fitted using n training data);
- 4 Define the N -dimensional vector \mathbf{Y} formed by the combination of the n -dimensional vector \mathbf{Y}^{num} (of n numerically simulated H_S data) and \tilde{Y}^{meta} . Given the GPD threshold t , construct \mathbf{Y}_t the vector of excesses $Y|Y>t$;

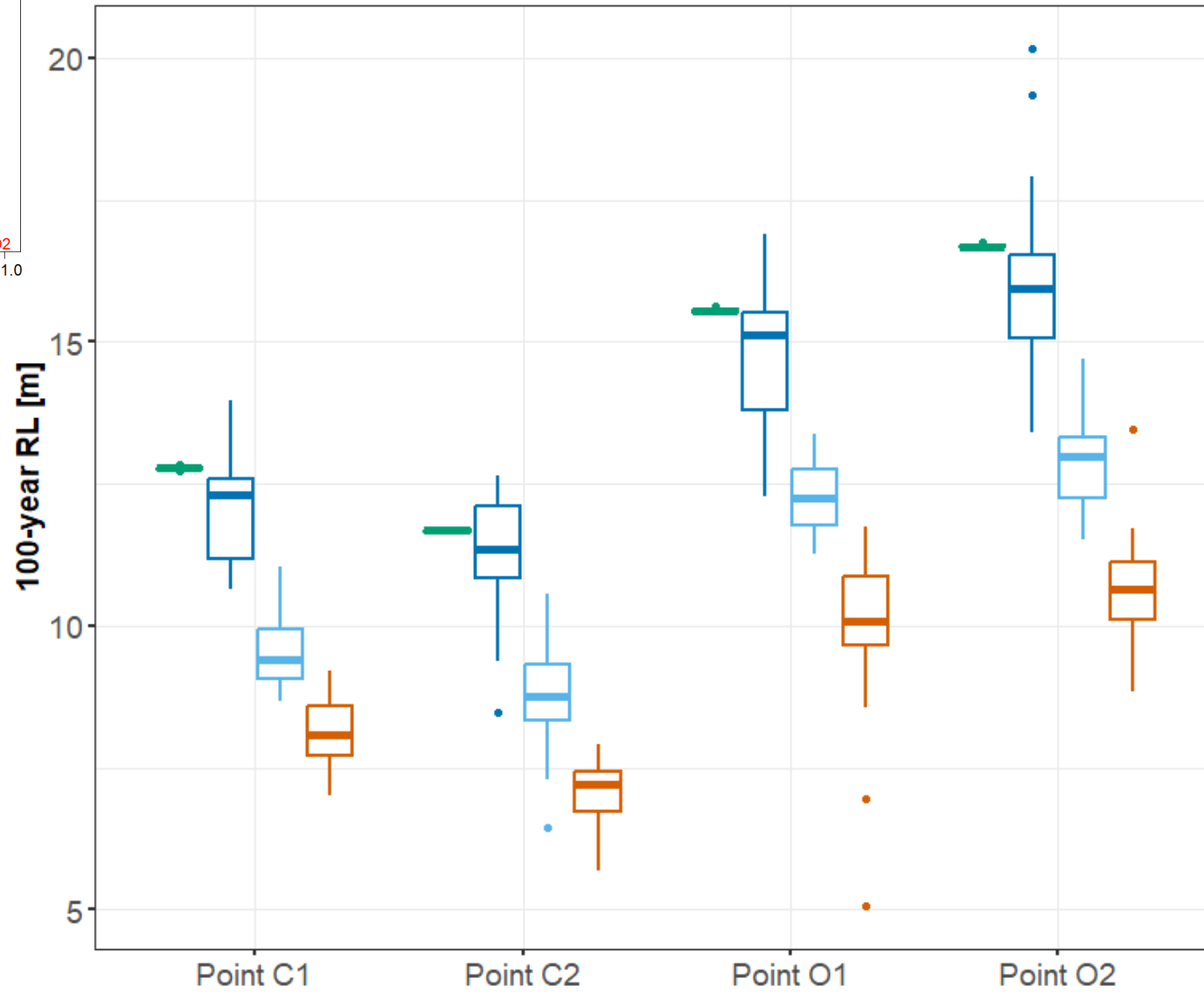
- 5 Simulate a N -dimensional vector from the GPD model given θ , $Y(\theta) \sim f(Y|\theta)$;
- 6 Accept θ if and only if $W(\mathbf{Y}_t, Y(\theta)) \leq \varepsilon$.

Acceptance / Rejection algorithm based on the Wasserstein distance

Probabilistic ML approach RFwE



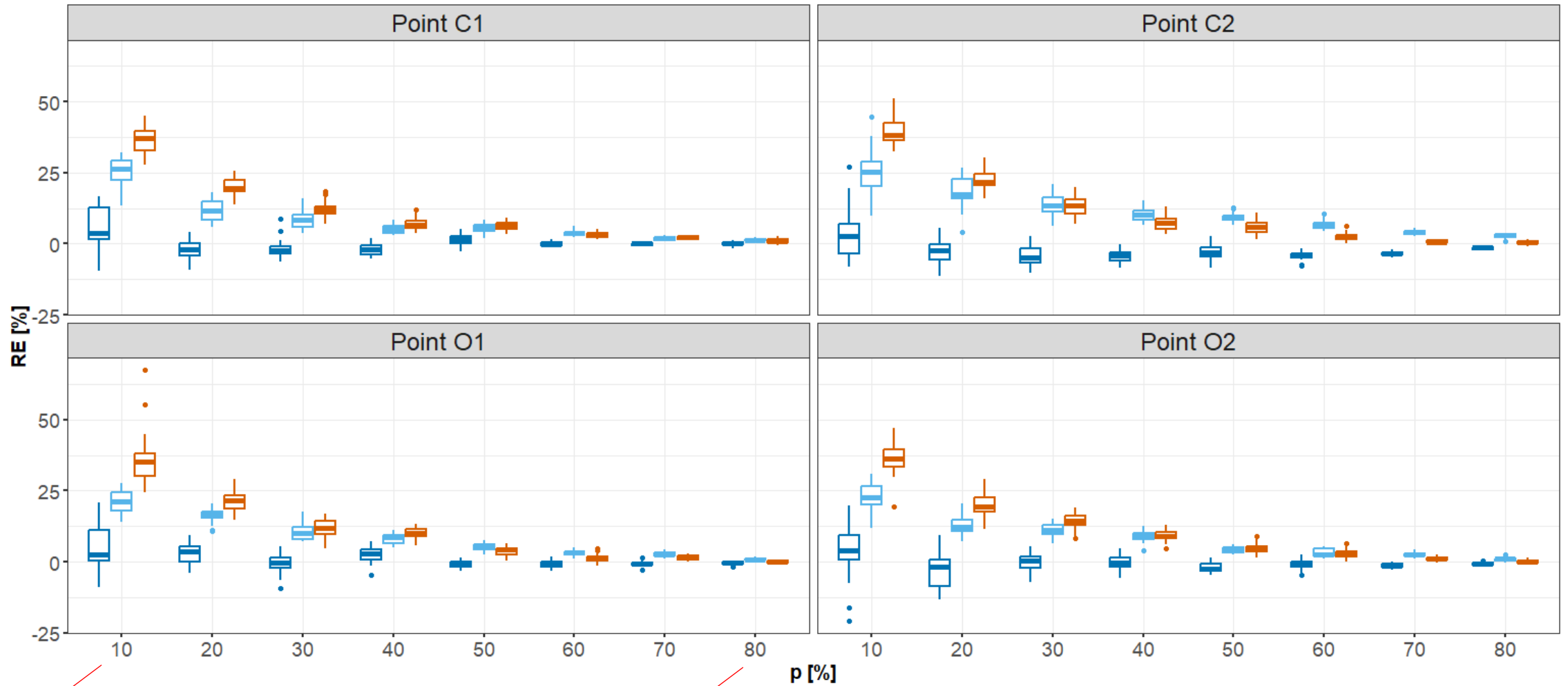
reference All RFwE RFwoE Subset



Relative error RE w.r.t. the proportion p of samples / full dataset

25 répétitions of the experiment

RFwE RFwoE Subset



70 cyclones

560 cyclones

16

Summary and open questions

- Estimates of 100 year Return Level **imposes the use of a large dataset (500 – 1,000)** of computed Hs maps; each of them corresponding to a synthetic Tropical Cyclone
- To perform the analysis using 50-100 Hs maps, we propose to augment the available database with **probabilistic ML predictions with account for ML prediction error**
- **Results suggest that the bias is largely lowered as well as the uncertainty estimates (95% confidence interval width)**
- **Open questions**
 - Could a limited number of synthetic tropical cyclones be **selected beforehand?**
 - Is it worth integrating the uncertainty related to the **hyperparameter tuning?**

Thank for your attention!

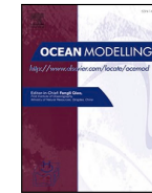
Ocean Modelling 186 (2023) 102275



Contents lists available at [ScienceDirect](#)

Ocean Modelling

journal homepage: www.elsevier.com/locate/ocemod



More reading?

Combining uncertain machine learning predictions and numerical simulation results for the extreme value analysis of cyclone-induced wave heights – Application in Guadeloupe

Jeremy Rohmer^{*}, Andrea G. Filippini, Rodrigo Pedreros

