

# **Large Scale Probabilistic and Multi-output Benthic Habitat Mapping**

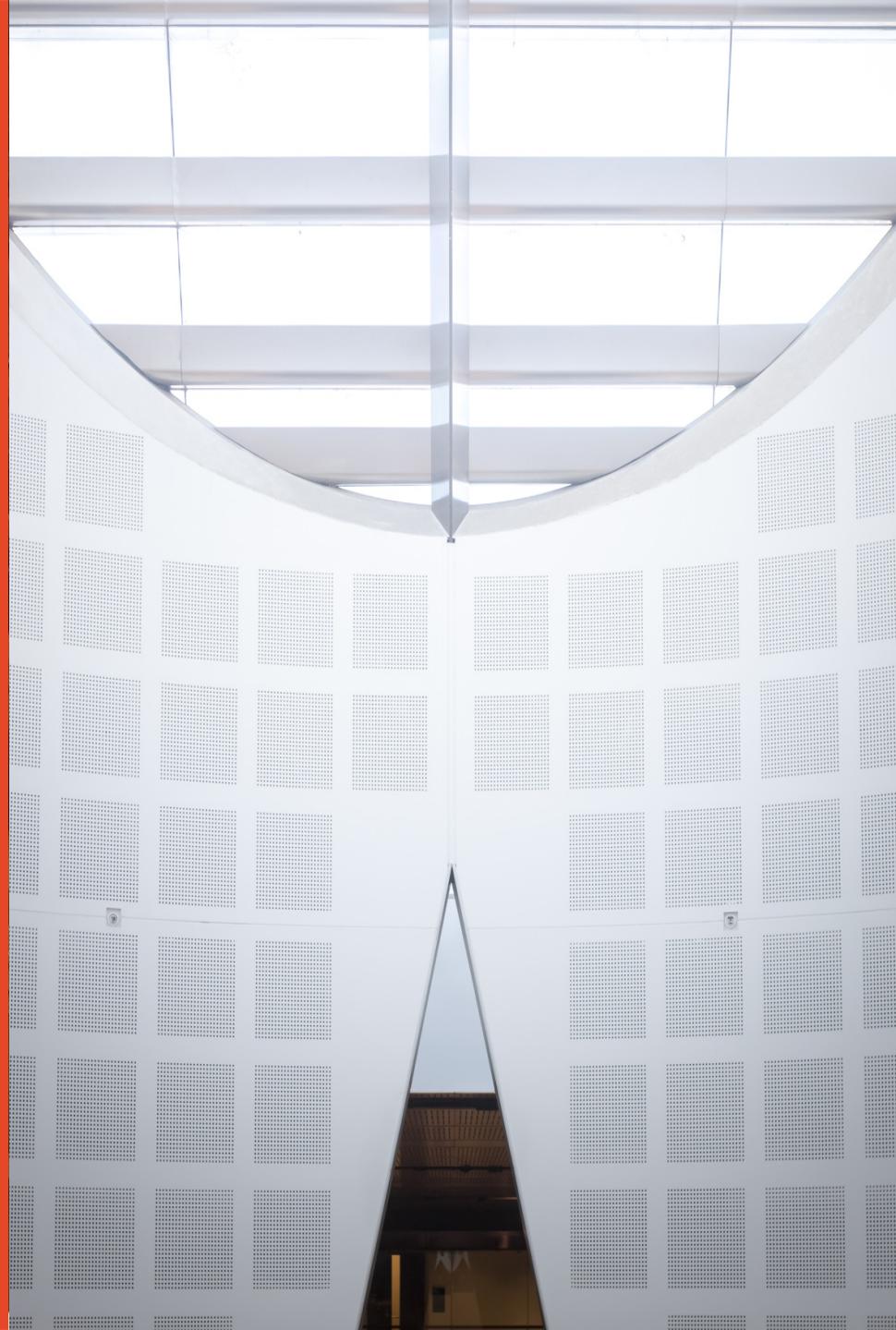
Honours Thesis

## **Presented by**

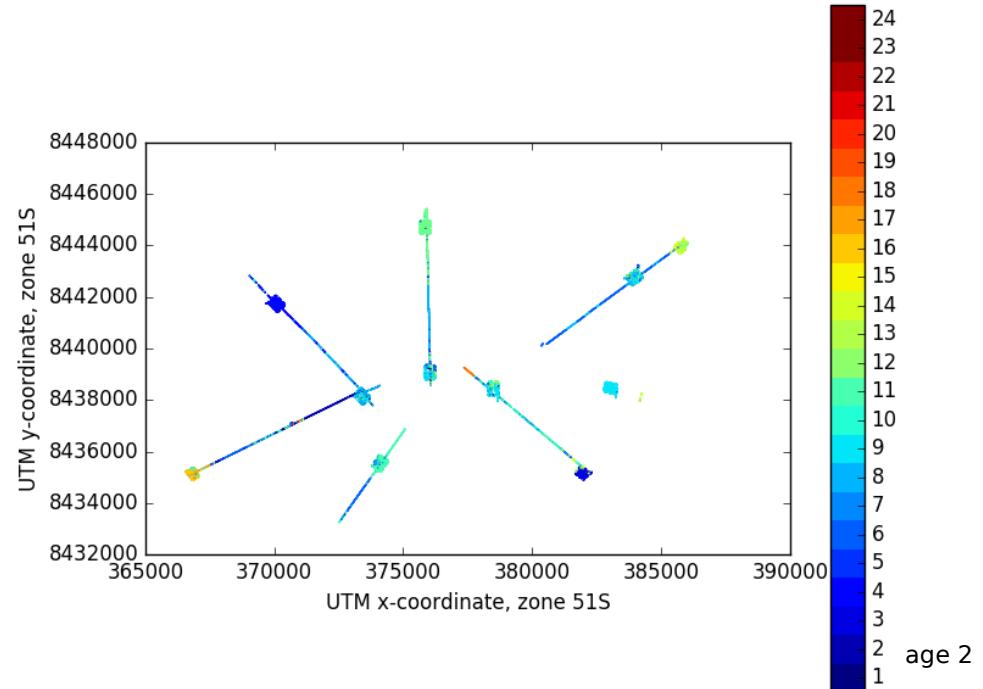
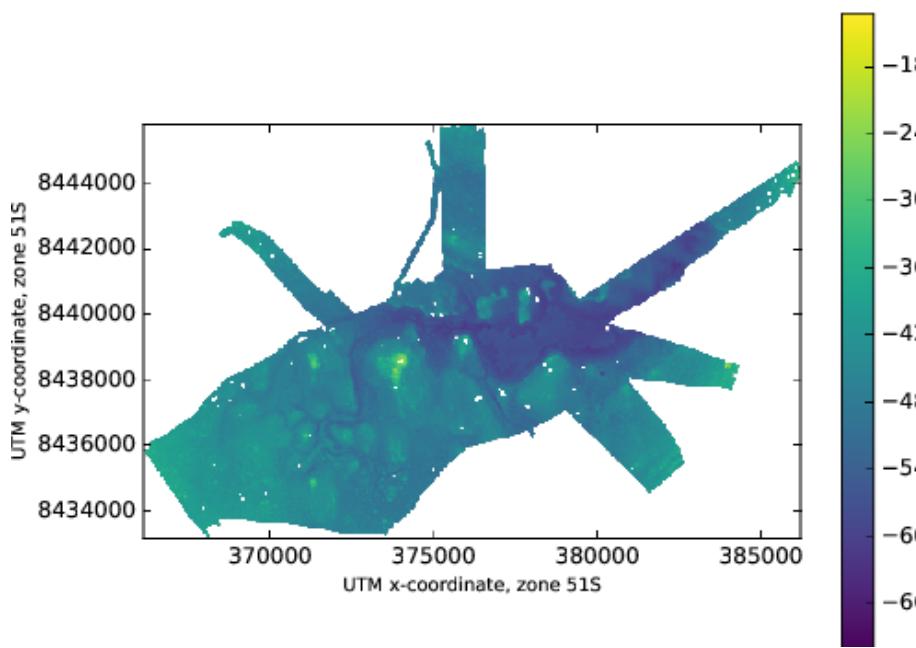
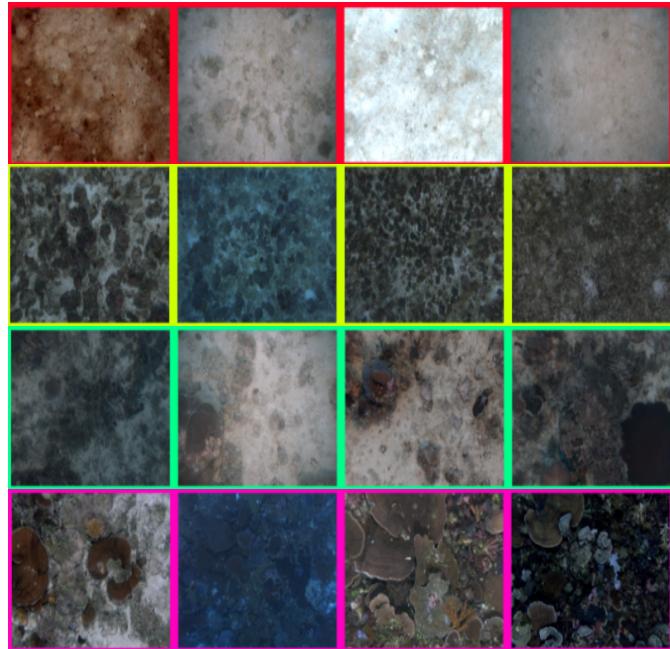
Justin Ting, School of Engineering and IT  
Supervisor: Simon O'Callaghan



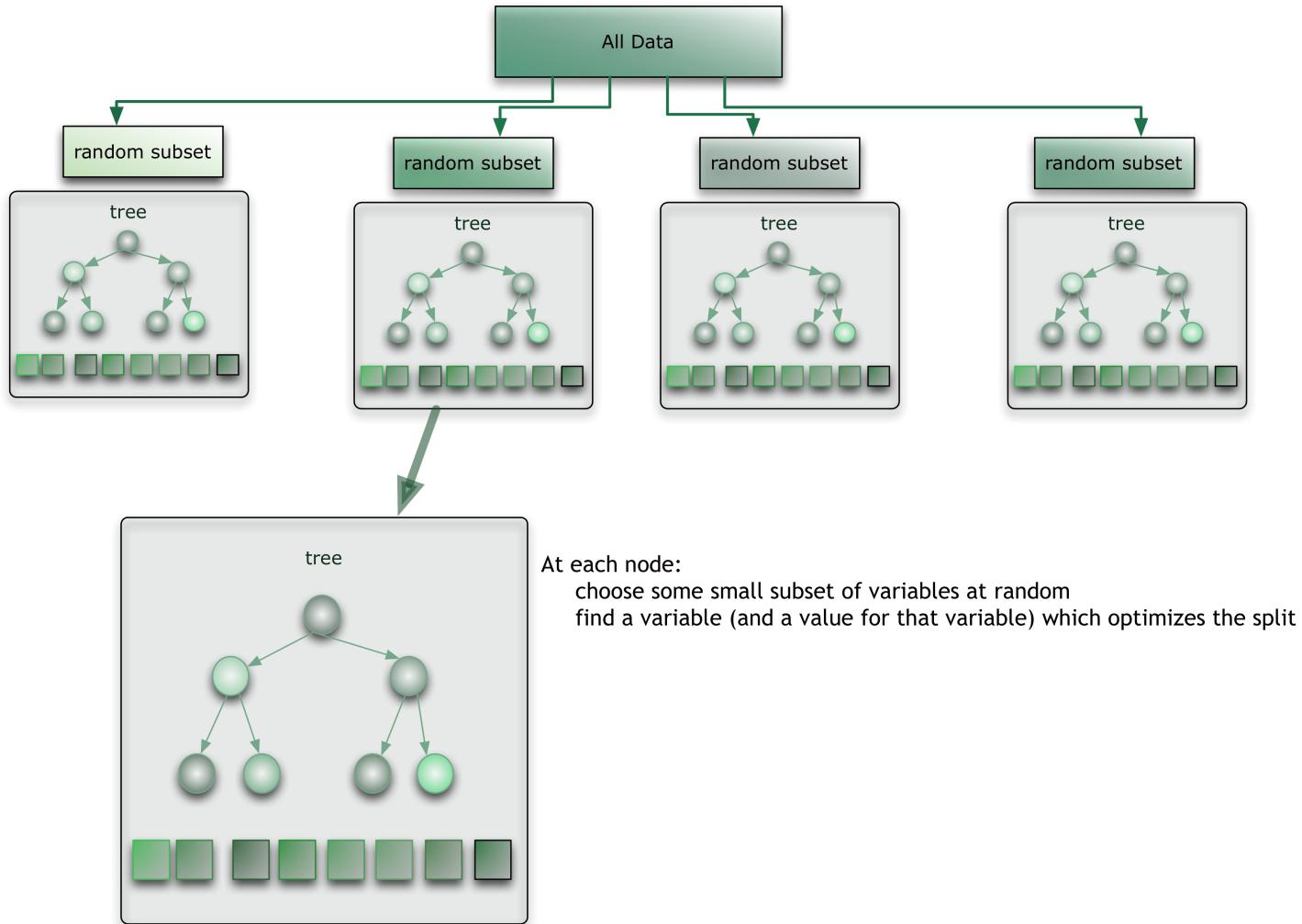
THE UNIVERSITY OF  
**SYDNEY**



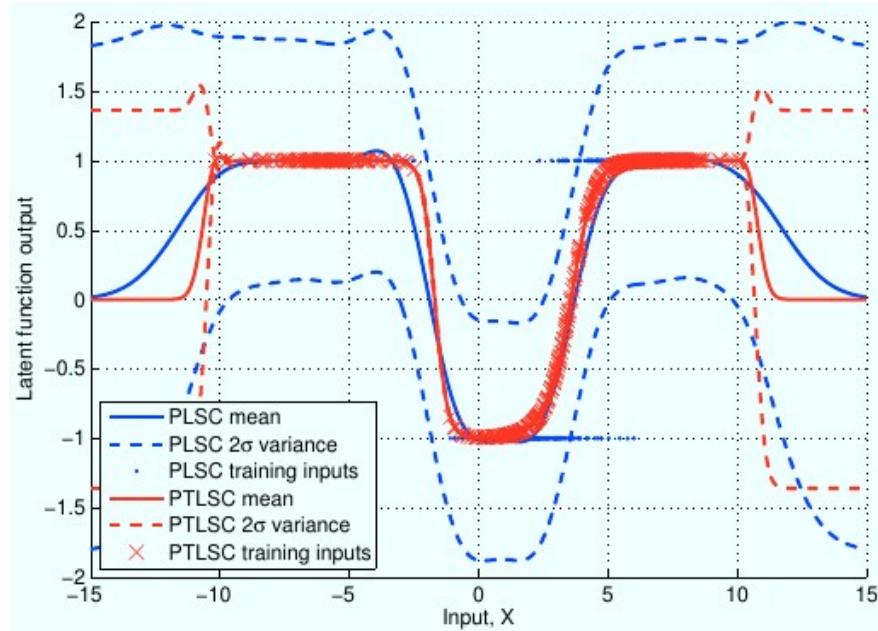
# Motivation



# Related Literature - Deterministic Machine Learning Methods



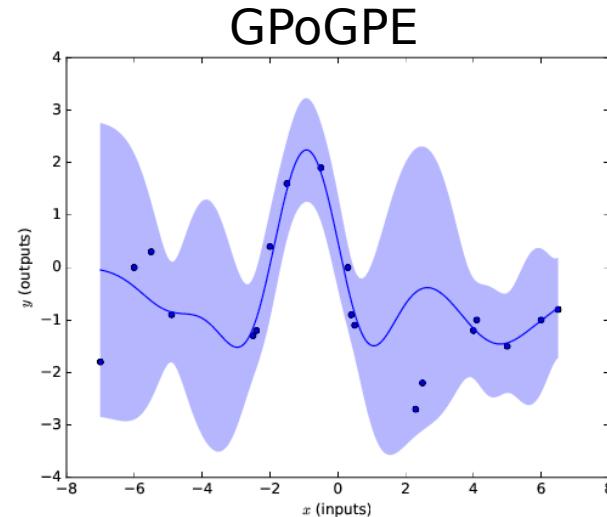
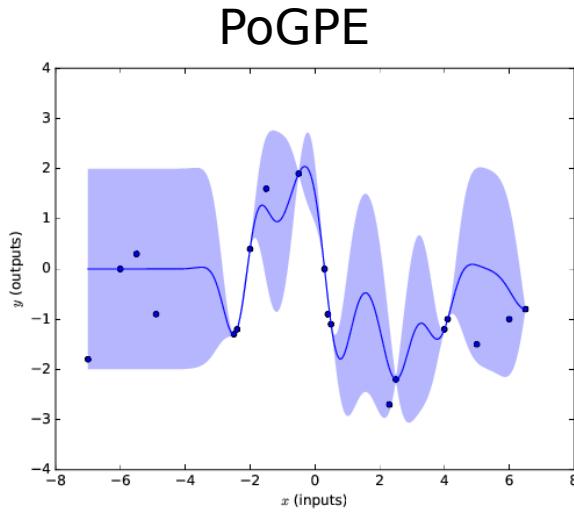
# Related Literature - Probabilistic Machine Learning Methods



$$\mathbf{K}_{\mathbf{x}_p, \mathbf{x}_q} = \sigma_f^2 e^{-\frac{1}{2}(\mathbf{x}_p - \mathbf{x}_q)^T L^{-1} (\mathbf{x}_p - \mathbf{x}_q)} + \sigma_n^2 \mathbf{I}$$

$$\mu_i = \mathbf{y}_i - \frac{[\mathbf{K}^{-1} \mathbf{y}]_{-i}}{[\mathbf{K}^{-1}]_{-ii}} \quad \text{and} \quad \sigma_i^2 = \frac{1}{[\mathbf{K}^{-1}]_{-ii}}$$

# Approach - Gaussian Process Approximation Methods



$$O(n_1^3) + O(n_2^3) + \dots + O(n_{k-1}^3) + O(n_k^3) \ll O(N^3)$$

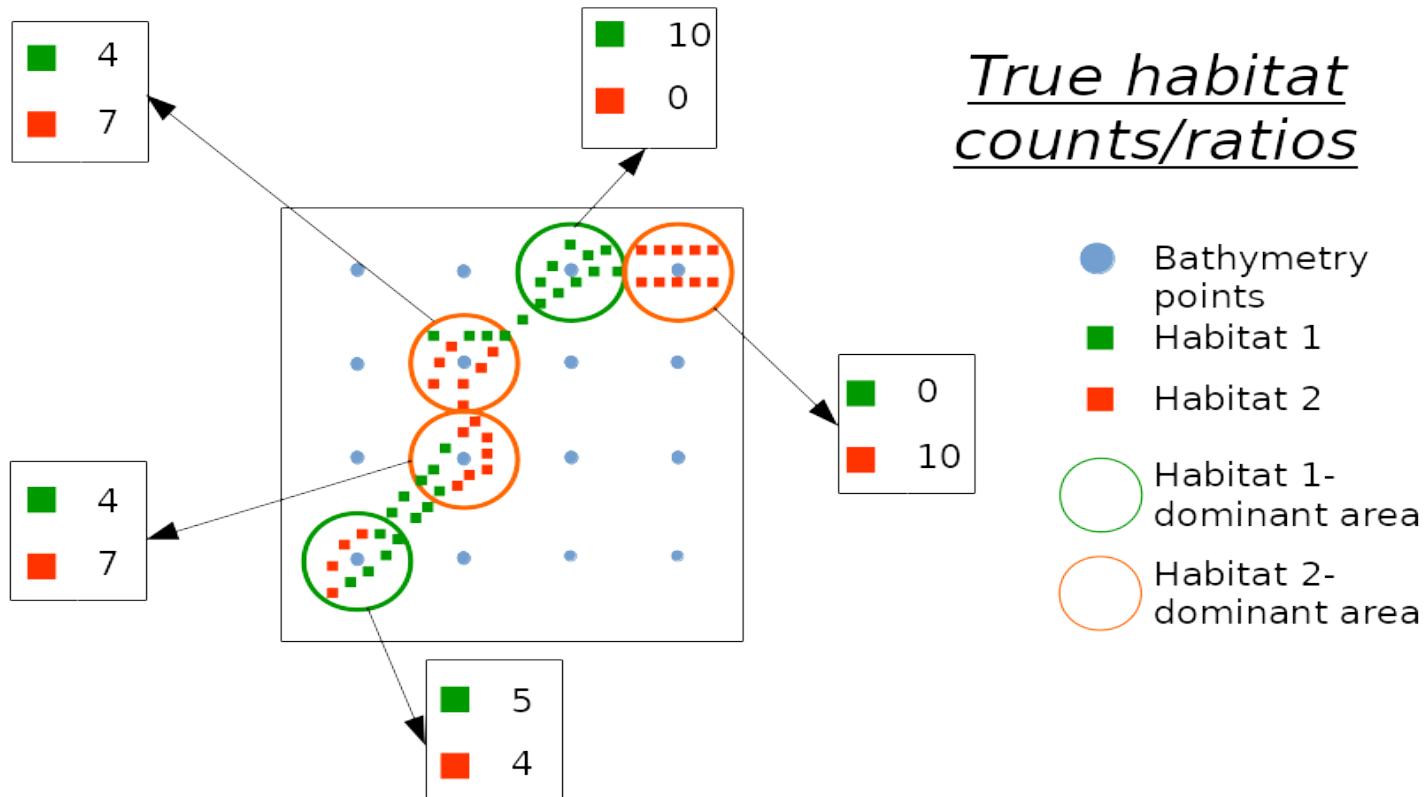
where  $n_1 + n_2 + \dots + n_{k-1} + n_k = N$

if  $n_k = n_{k-1} = \dots = n_2 = n_1$ , and  $k = 100, N = 20000$

$$100 \times O(200^3) \ll O(20000^3) \rightarrow 8 \times 10^8 \ll 8 \times 10^{12}$$

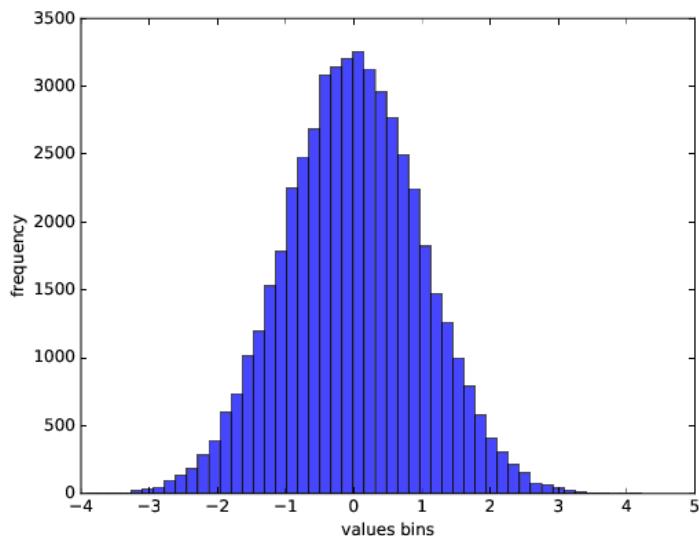
# Approach - Intro to Multi-output Methods

- However, the data collection process does not often align exactly with this

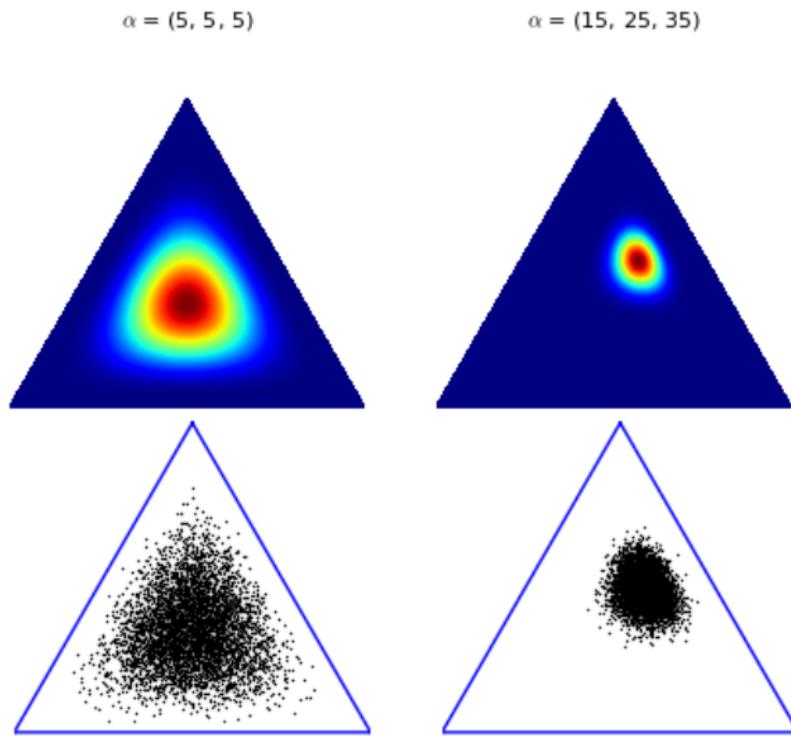


# Approach - Dirichlet Multinomial Regression

Multinomial Distribution



Dirichlet Distribution



<todo - equation/s?>

# Results - GP vs DM vs Deterministic Methods

	SVM	Random Forest	Logistic Regression	K-Nearest Neighbour	Gaussian Process Classification	Generalised Product of Gaussian Process Experts	Dirichlet Multinomial
Aggregated (4) labels	Accuracy	<u>0.75</u>	0.79	0.77	0.78	<b>0.83</b>	0.75
	F-score	<u>0.21</u>	0.47	0.34	0.47	<b>0.53</b>	0.33
Full (24) labels	Accuracy	0.29	0.35	0.31	0.33	<b>0.39</b>	<u>0.21</u>
	F-score	<u>0.10</u>	0.22	0.13	0.26	<b>0.32</b>	0.19

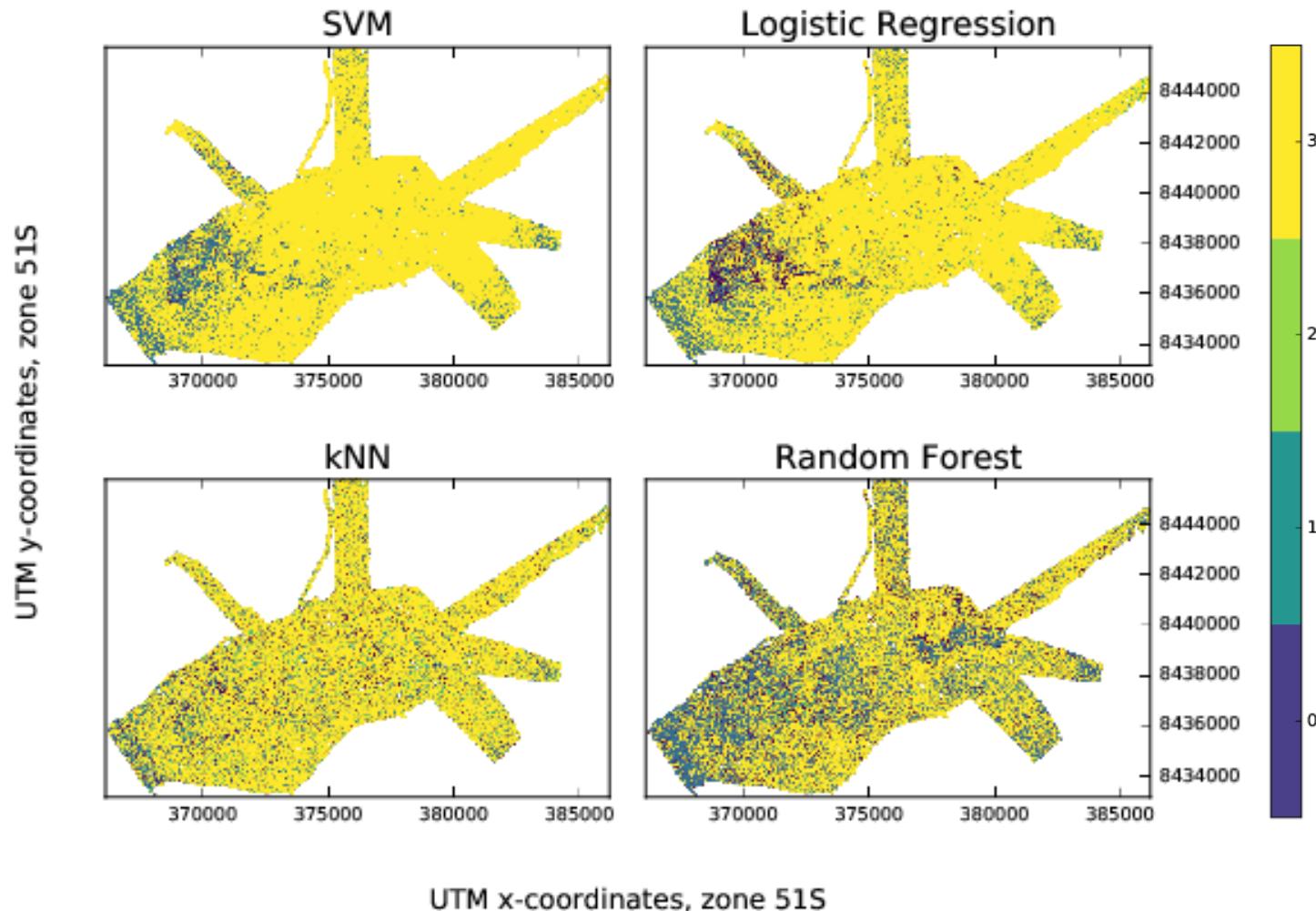
# Results

## Time of GP vs approx GP vs DM

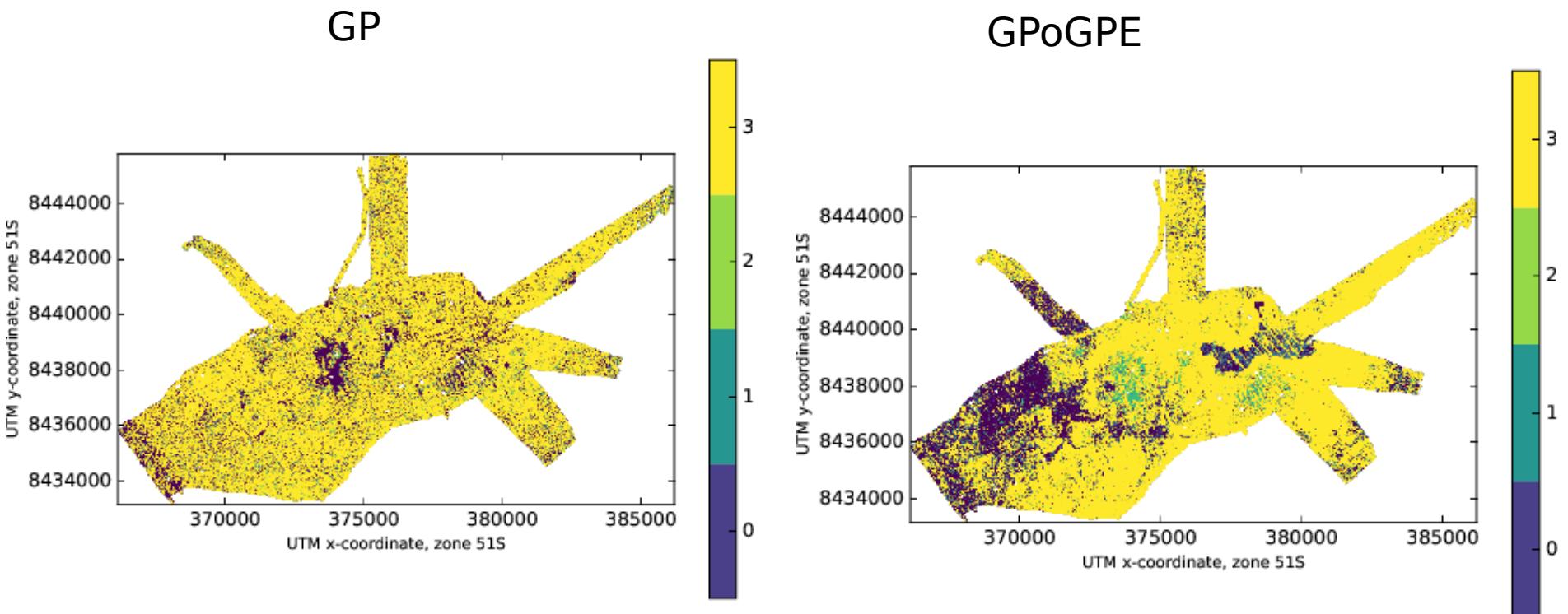
Time taken to perform model fitting and predictions on 4-habitats, with 4700 training points, and 500,000 test points

<b>Heading 1</b>	<b>Gaussian Process Classification</b>	<b>Generalised Product of Gaussian Process Experts</b>	<b>Dirichlet Multinomial Regression</b>
Model fitting	47:41	00:23	00:01
Predictions	20:41	03:35	00:05

# Discussion - Deterministic Maps

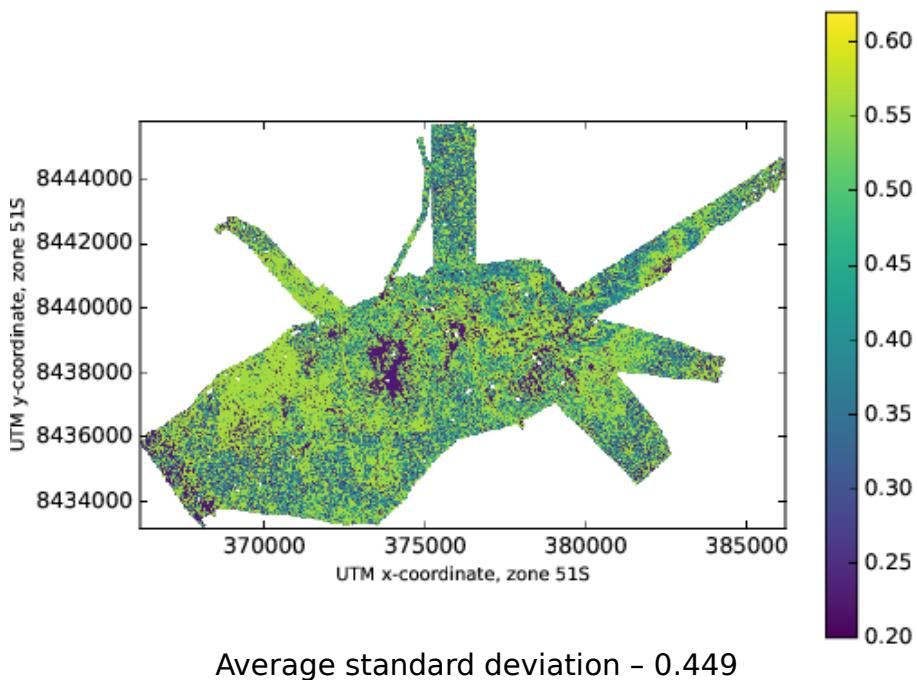


# Discussion - Gaussian Process and Approximations

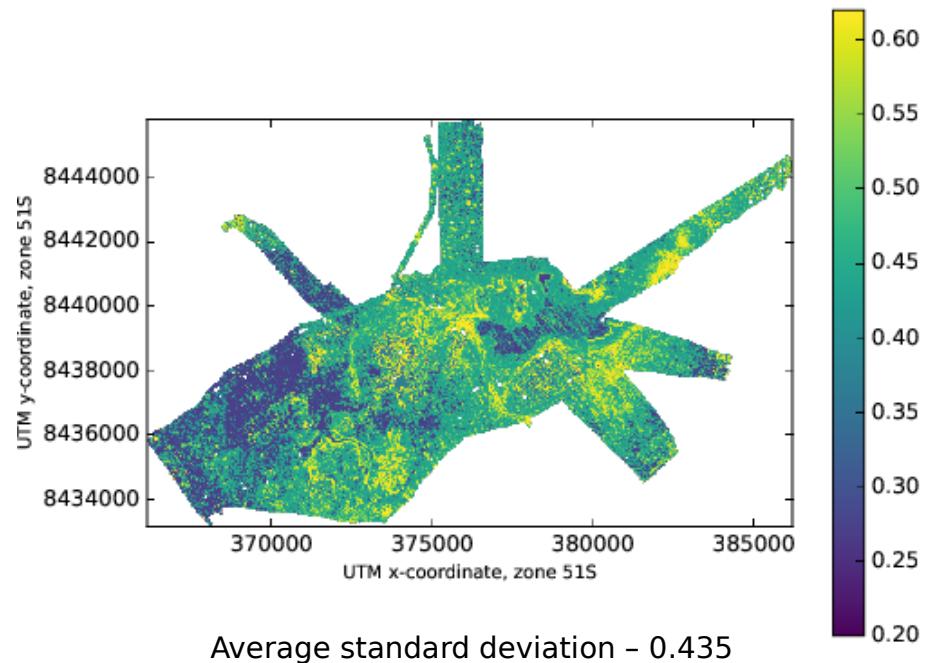


# Discussion - Gaussian Process Variance

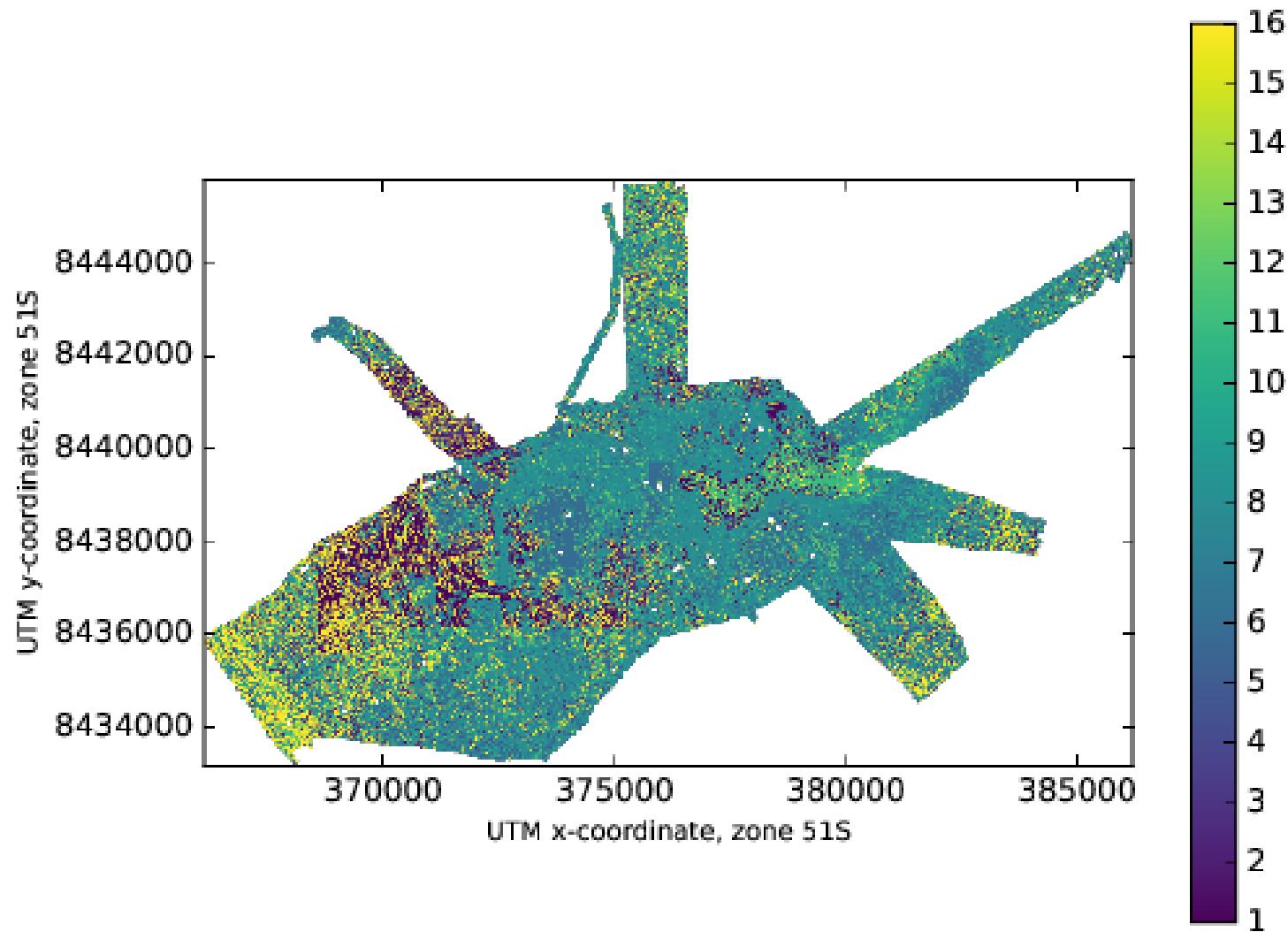
Gaussian Process Standard Deviation at most likely labels



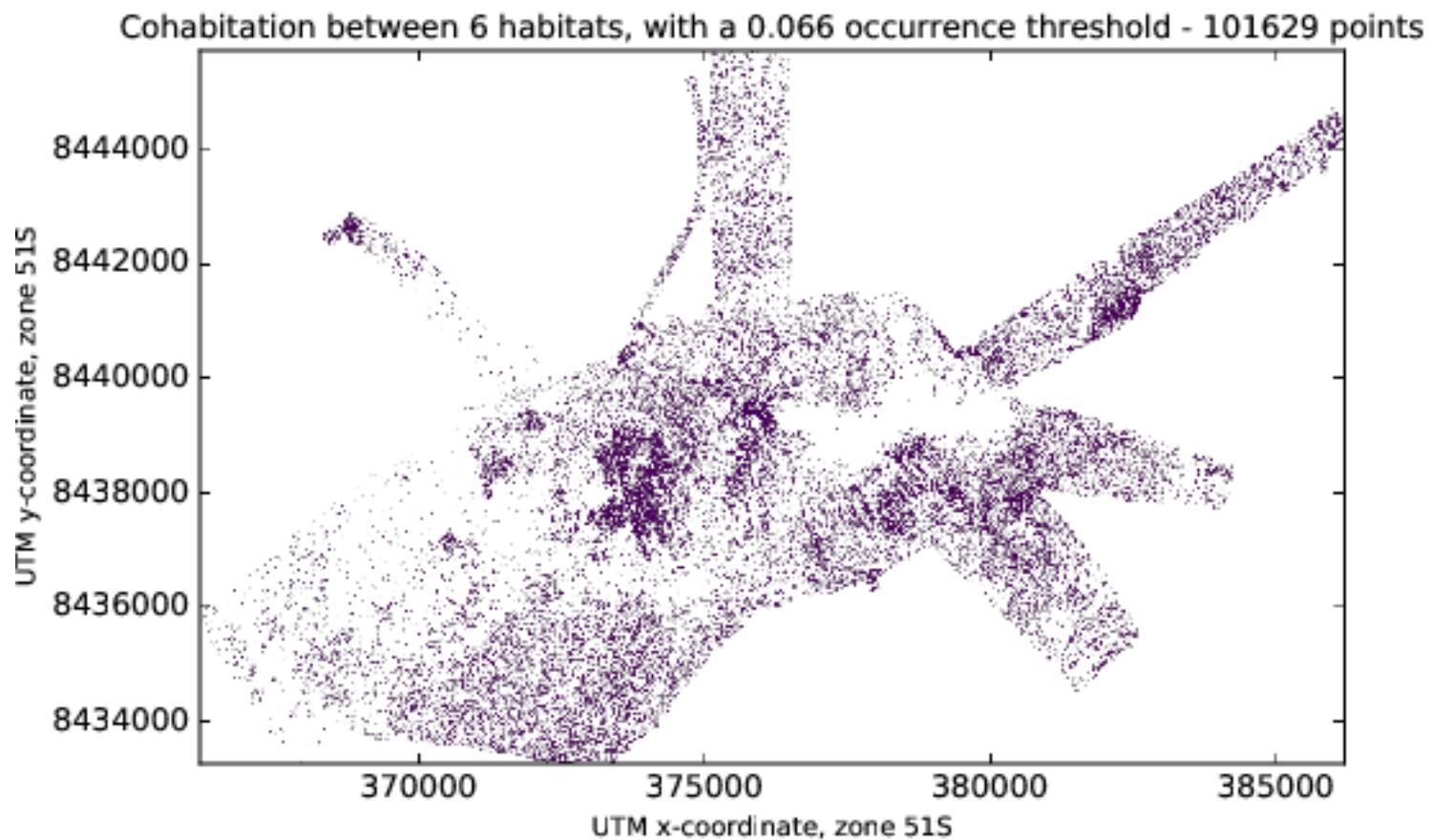
Generalised Product of Gaussian Process Experts Standard Deviation at most likely labels



# Discussion – Dirichlet Multinomial Regression



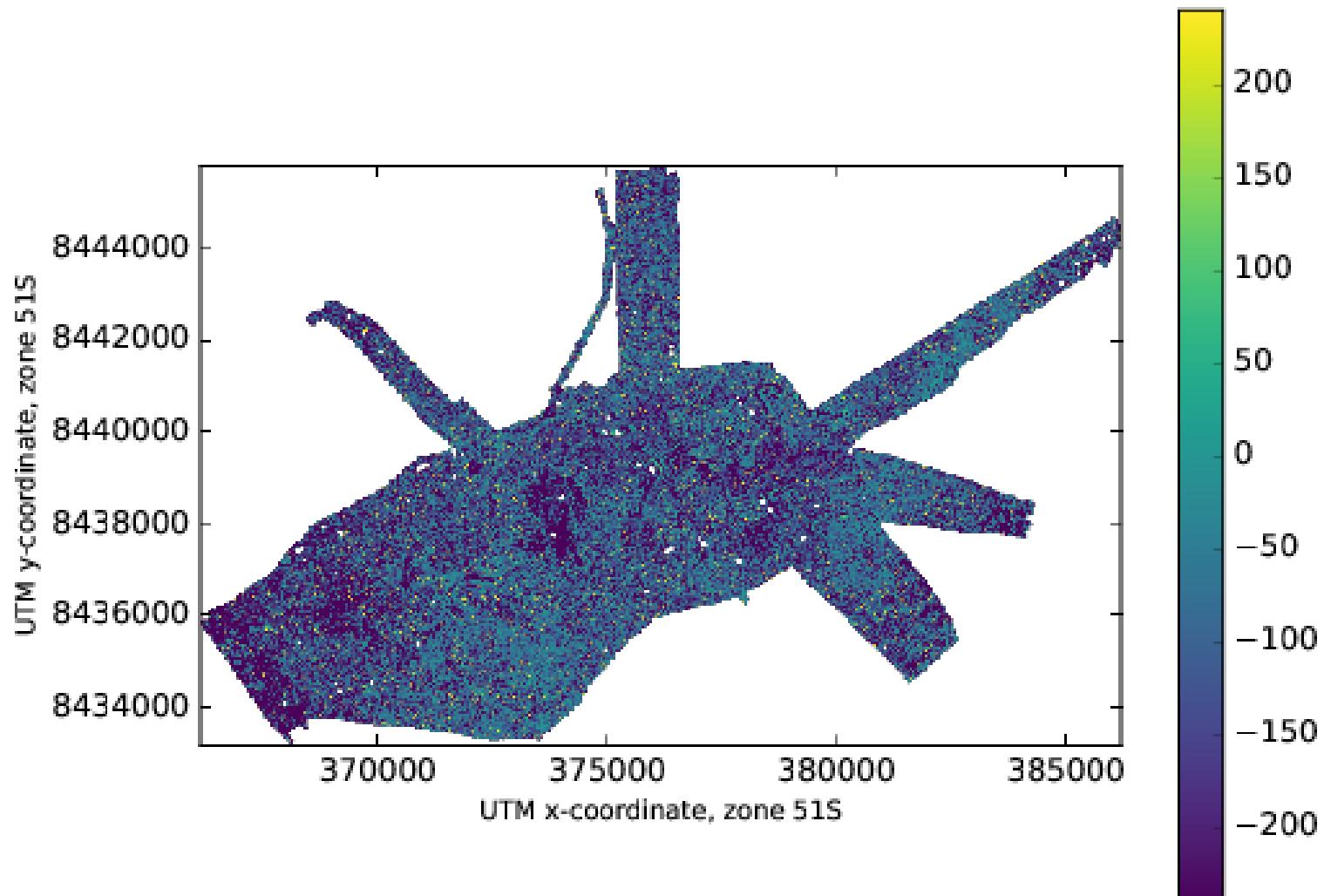
# Discussion - Biodiversity



# Discussion - Dirichlet Multinomial using parameters from full distribution

<todo - argmax gif>

# Discussion - Dirichlet Multinomial Entropy



# Limitations

- Marine biologist/expert verification of map quality
- 
-

# Conclusion

- Lorem Ipsum