

# Ecosystem Services

## Reducing Uncertainty in Ecosystem Service Modelling through Weighted Ensembles

--Manuscript Draft--

<b>Manuscript Number:</b>	ECOSER-D-21-00429R2
<b>Article Type:</b>	Research Paper
<b>Keywords:</b>	Committee averaging; Prediction Error; Accuracy; United Kingdom; Validation; Weighted averaging
<b>Corresponding Author:</b>	Simon Willcock Bangor University Bangor, Gwynedd UNITED KINGDOM
<b>First Author:</b>	Danny A.P. Hooftman
<b>Order of Authors:</b>	Danny A.P. Hooftman James M. Bullock Laurence Jones Felix Eigenbrod José I. Barredo Matthew Forrest Georg Kindermann Amy Thomas Simon Willcock
<b>Abstract:</b>	<p>Over the last decade many ecosystem service (ES) models have been developed to inform sustainable land and water use planning. However, uncertainty in the predictions of any single model in any specific situation can undermine their utility for decision-making. One solution is creating ensemble predictions, which potentially increase accuracy, but how best to create ES ensembles to reduce uncertainty is unknown and untested. Using ten models for carbon storage and nine for water supply, we tested a series of ensemble approaches against measured validation data in the UK. Ensembles had at minimum a 5-17% higher accuracy than a randomly selected individual model and, in general, ensembles weighted for among model consensus provided better predictions than unweighted ensembles. To support robust decision-making for sustainable development and reducing uncertainty around these decisions, our analysis suggests various ensemble methods should be applied depending on data quality, for example if validation data are available.</p>
<b>Suggested Reviewers:</b>	Becky Chaplin-Kramer bchaplin@stanford.edu ES modelling expert  Rachel Neugarten rachel.neugarten@gmail.com ES expert  Javier Martinez-Lopez jmartinez@cebas.csic.es ES modelling expert  Benjamin Bryant bpbryant@stanford.edu Expert in ES modelling & uncertainty  Nynke Schulp nynke.schulp@vu.nl ES expert
<b>Response to Reviewers:</b>	

**Declaration of interests**

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:



PRIFYSGOL  
**BANGOR**  
UNIVERSITY

2 December 2021

Dear Editor,

***Reducing Uncertainty in Ecosystem Service Modelling through Weighted Ensembles***

I am pleased submit the above paper for your consideration, which has been revised in accordance with comments from both the Editors and two reviewers. This manuscript presents a study, unprecedented in scope, on maximising the accuracy of ensembles of models of ecosystem services. As such, we feel our paper is a good fit for Ecosystem Services given that the journal's content seeks to understand science, policy and practice of Ecosystem Services.

The important knowledge-gap we address in our manuscript is that global efforts to quantify ecosystem services (e.g. through the Intergovernmental Science-Policy Platform on Biodiversity and Ecosystem Services [IPBES]) are lagging behind those of other grand challenges (e.g. the Intergovernmental Panel on Climate Change [IPCC]). For example, whilst the IPCC use ensembles of models to provide robust estimates of plausible futures, the latest state-of-the-art ES models produced via IPBES rely on single model outputs with little/no validation (e.g., see [Chaplin-Kramer et al., 2019](#)). This is because, unlike climate models, ecosystem service models often differ in the forms of their outputs – even when modelling the same services. As a result, it is currently not known how best to combine distinct ecosystem service model outputs to provide reliable ensemble products. In this manuscript we show how best to overcome these issues.

Our study – which uses ten models for carbon storage and nine for water supply, to test ten contrasting ensemble approaches against 2,597 validation data points in the UK – is the first assessment of different approaches to creating ecosystem service model ensembles. Our findings represent important advances of significance to scientists and policy-makers working within ecosystem services, environmental science and sustainability, as well as the wider natural science modelling community.

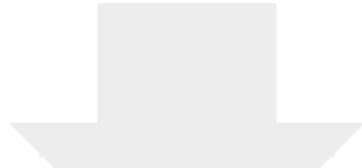
We show that using an individual ecosystem service model is fraught with concerns as *a priori* it is not known which is the most accurate and choosing only one model can, at worst, result in perverse decisions. Deriving decisions from an ensemble of ES models provides an improvement over using one model for any location, but also more consistency over larger scales. Using weighted average ensemble approaches further improves accuracy but also substantially decreases uncertainty among ensemble approaches compared to uncertainty among models, a further indication of increased fit to reality. Thus, particularly when validation data are not available, we recommend the use of weighted ensembles in ES research to substantially reduce uncertainty and to support robust decision-making for sustainable development.

In partnership with decision-makers, the important advances suggested in our manuscript could help to ensure ecosystem service research contributes to and informs ongoing policy processes (such as IPBES, the Sustainable Development Goals and CBD Aichi targets) and facilitates the development of indicators for the monitoring of human well-being in United Kingdom and beyond.

We thank you in advance for considering our manuscript.

Yours sincerely,

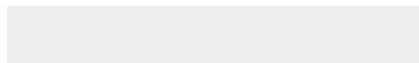
Prof Simon Willcock



Click here to access/download

**e-Component**

WeightedEnsembles\_SI\_Rev\_v4\_CLEAN.docx



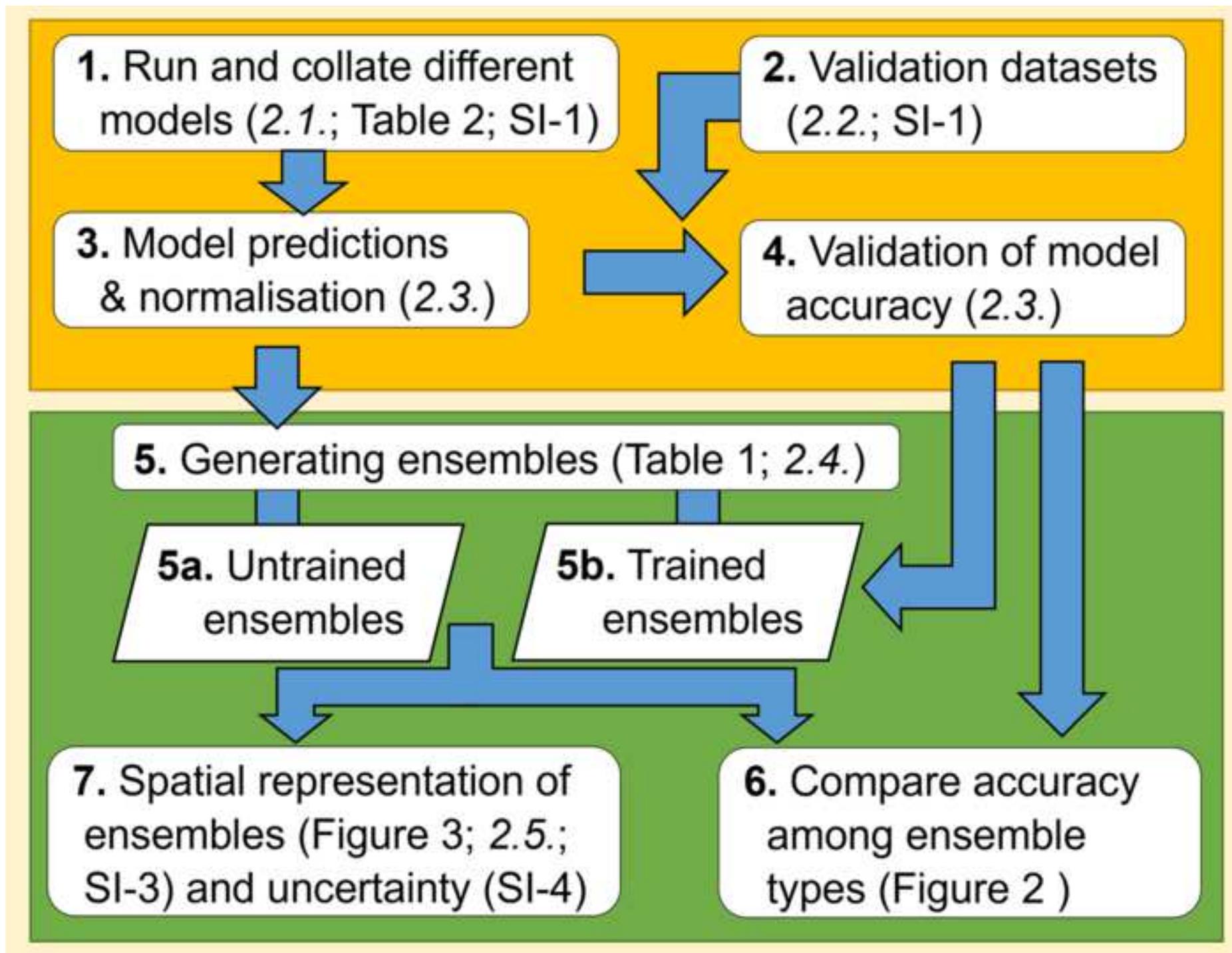


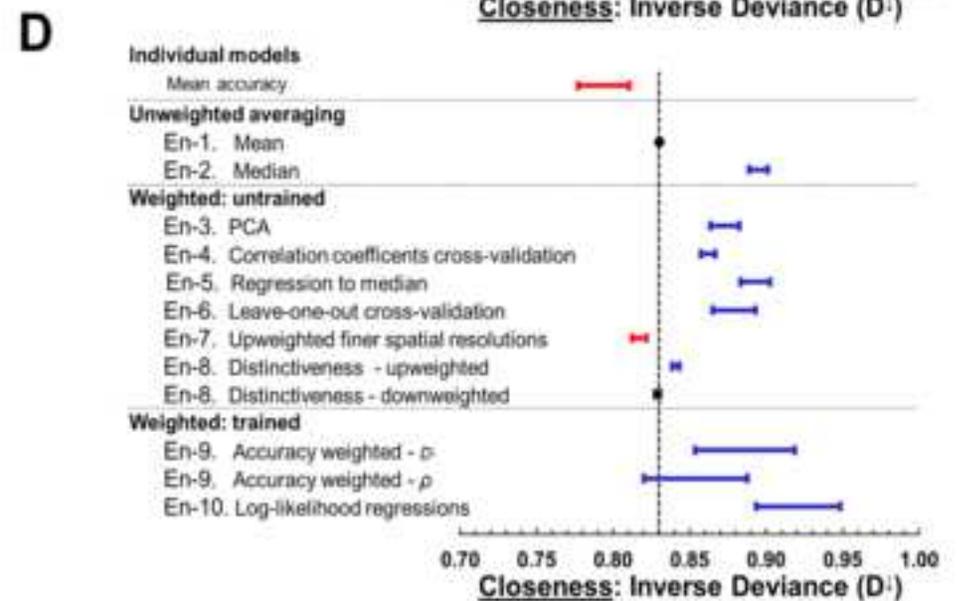
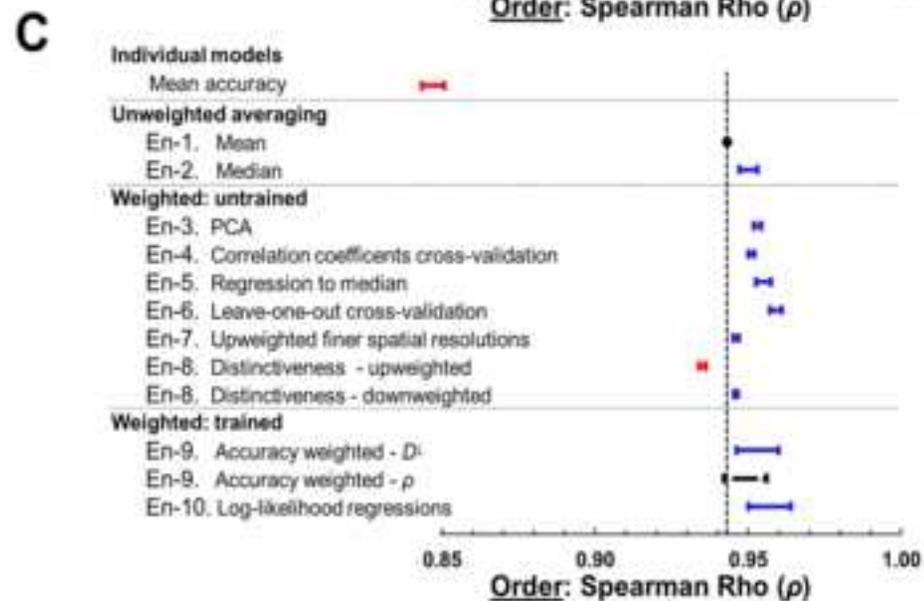
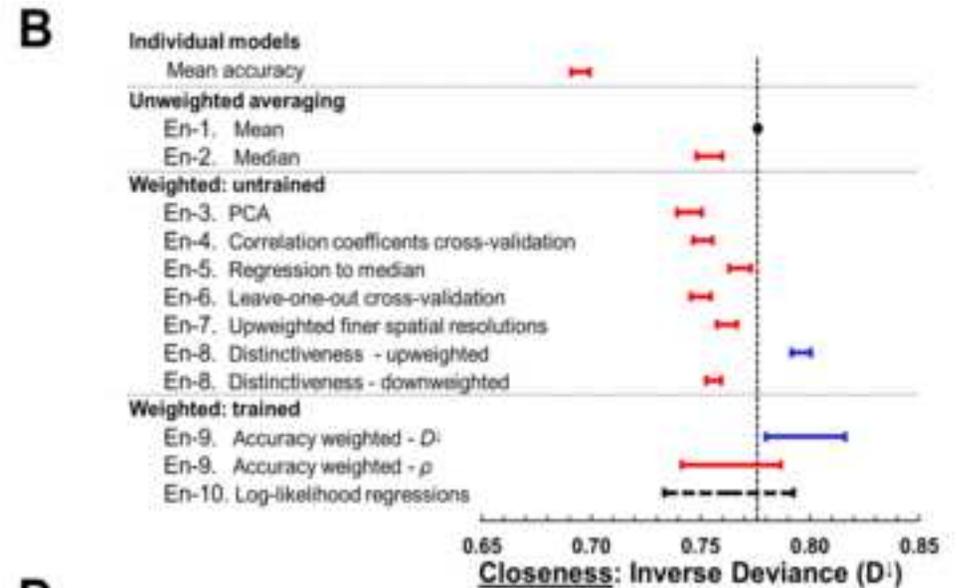
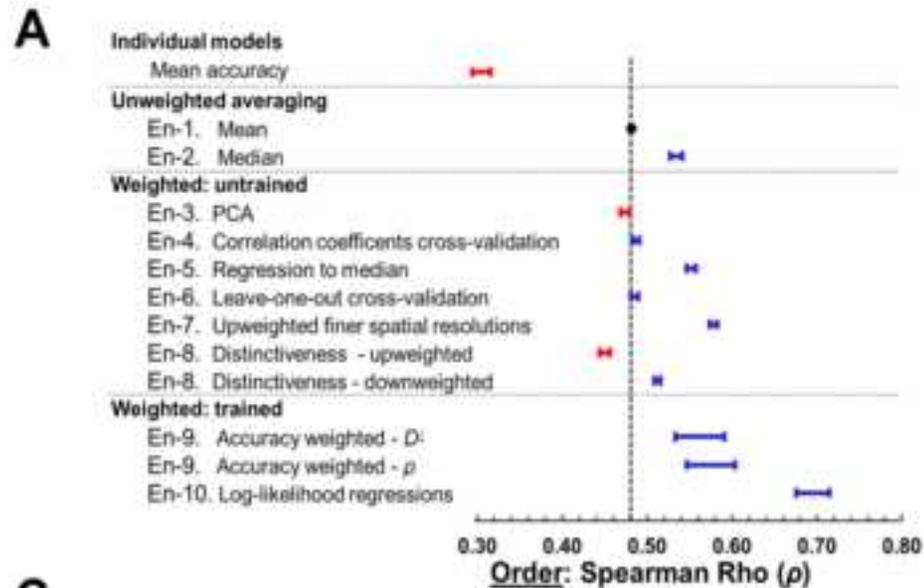
Click here to access/download

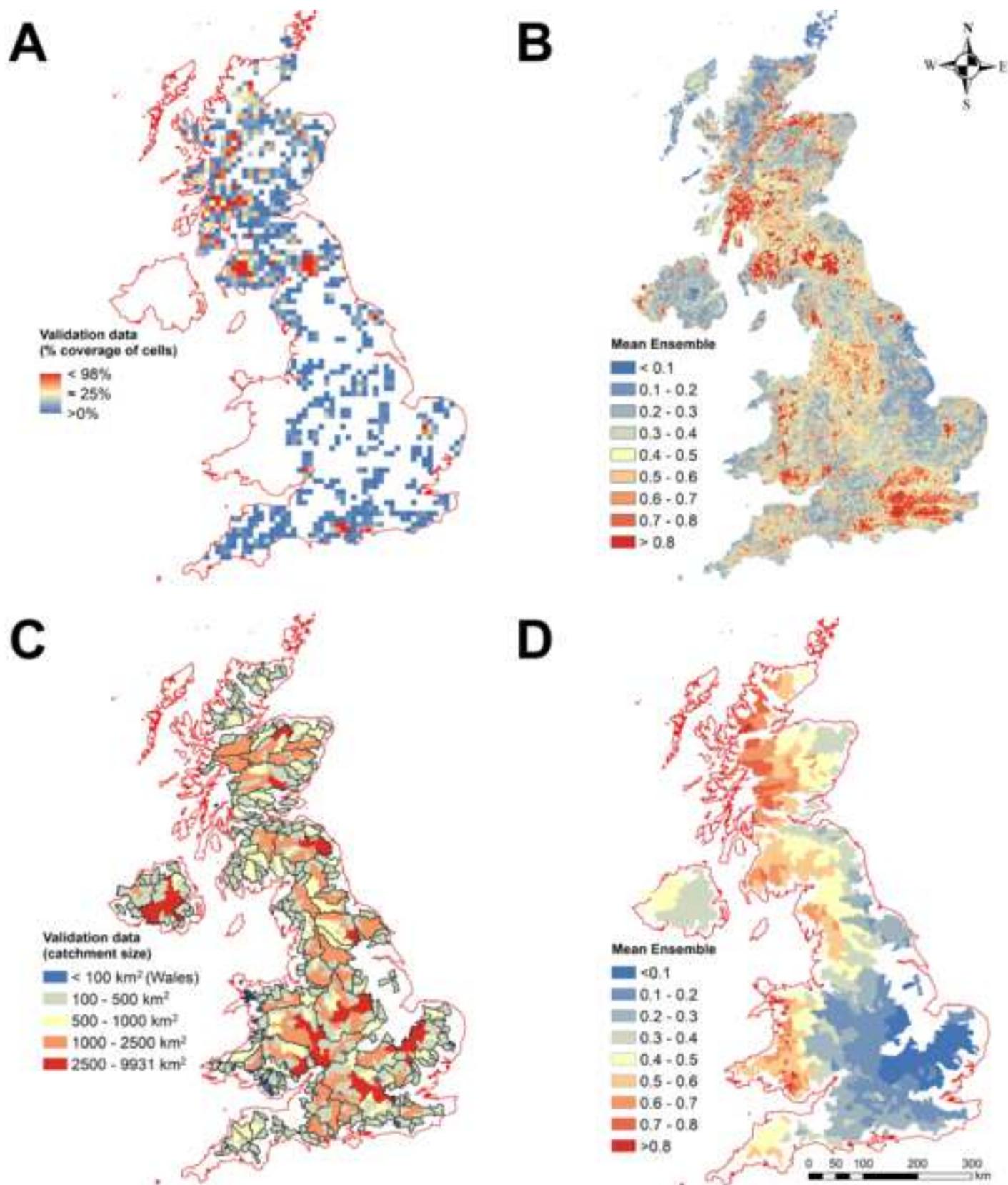
**e-Component**

WeightedEnsembles\_SI\_Rev\_v4.docx









## Accuracy compared to mean ensemble

**Individual Models**

**Mean accuracy is always worse**

**Unweighted Ensembles**

Mean Ensemble

Reference ensemble type: up to 19% better than individual models

Median Ensemble

**Mostly better, rarely worse**

**Untrained Weighted Ensembles**

Deterministic Consensus

**Sometimes better, sometimes worse**

Iterated Consensus

**Mostly better, rarely worse**

Attribute based

**Sometimes better, sometimes worse**

**Trained Weighted Ensembles**

Accuracy weighted

**Mostly better, rarely worse**

Regressed consensus

**Mostly better, never worse**

**Highlights:**

- Ensembles of models are used for other disciplines but not ecosystem services
- How best to combine ecosystem service models into an ensemble is unknown
- We test ten contrasting ensemble approaches
- Ensembles had up to 27% higher accuracy than a randomly selected individual model
- Weighted ensembles provided better predictions

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30

# ANONYMISED MANUSCRIPT

## Reducing Uncertainty in Ecosystem Service Modelling through Weighted Ensembles

### Highlights:

- Ensembles of models are used for other disciplines but not ecosystem services
- How best to combine ecosystem service models into an ensemble is unknown
- We test ten contrasting ensemble approaches
- Ensembles had up to 27% higher accuracy than a randomly selected individual model
- Weighted ensembles provided better predictions

### Abstract: (150 words)

Over the last decade many ecosystem service (ES) models have been developed to inform sustainable land and water use planning. However, uncertainty in the predictions of any single model in any specific situation can undermine their utility for decision-making. One solution is creating ensemble predictions, which potentially increase accuracy, but how best to create ES ensembles to reduce uncertainty is unknown and untested. Using ten models for carbon storage and nine for water supply, we tested a series of ensemble approaches against measured validation data in the UK. Ensembles had at minimum a 5-17% higher accuracy than a randomly selected individual model and, in general, ensembles weighted for among model consensus provided better predictions than unweighted ensembles. To support robust decision-making for sustainable development and reducing uncertainty around these decisions, our analysis suggests various ensemble methods should be applied depending on data quality, for example if validation data are available.

### Graphical Abstract:

<b>Accuracy compared to mean ensemble</b>	
<b>Individual Models</b>	<b>Mean accuracy is always worse</b>
<b>Unweighted Ensembles</b>	
Mean Ensemble	Reference ensemble type: up to 19% better than individual models
Median Ensemble	<b>Mostly better, rarely worse</b>
<b>Untrained Weighted Ensembles</b>	
Deterministic Consensus	<b>Sometimes better, sometimes worse</b>
Iterated Consensus	<b>Mostly better, rarely worse</b>
Attribute based	<b>Sometimes better, sometimes worse</b>
<b>Trained Weighted Ensembles</b>	
Accuracy weighted	<b>Mostly better, rarely worse</b>
Regressed consensus	<b>Mostly better, never worse</b>

**Keywords:** Carbon; Committee averaging; Prediction Error; Accuracy; United Kingdom; Validation; Water supply; Weighted averaging

**Video Summary:** (see attached file)

## 31 1. Introduction

32 If the United Nations' sustainable development goals (SDG) are to be achieved worldwide (Griggs *et al.*  
33 2013), it is vital to understand and manage “*nature's contributions to people*” (termed ecosystem services;  
34 ES; Pascual *et al.* 2017). The empirical data needed to quantify ES are sparse in many parts of the world  
35 (Suich *et al.* 2015; Willcock *et al.* 2016), which is problematic as ES need to be accurately assessed and  
36 mapped to be incorporated in policy making and planning decisions (UKNEA 2011; de Groot *et al.* 2012).  
37 Such decisions require assessment of multiple ES, and the synergies and trade-offs among these ES, in order  
38 to estimate potential effects of land/water use change or other impacts (Willcock *et al.* 2016). Spatially-  
39 explicit models produce maps of estimated ES – typically based on globally available datasets of land cover  
40 combined with other predictor variables – and so can provide credible information of the spatial distributions  
41 of multiple ES, particularly where empirical data are lacking (Malinga *et al.* 2015; Costanza *et al.* 2017).  
42

43 Over the last 10 years, many ES models have been developed, by different teams, often using dissimilar  
44 approaches, and with little reference to the other models (Bagstad *et al.* 2013; Ochoa & Urbina-Cardona  
45 2017). For example, carbon stocks for climate change mitigation can be modelled by ‘look-up tables’  
46 relating land cover to stocks, by deterministic statistical inference, or by simulating complex processes  
47 (Willcock *et al.* 2019). However, most applications of ES models rely on only a single model for each ES  
48 (Englund *et al.* 2017; Bryant *et al.* 2018). Furthermore, while models can only approximate reality, few  
49 applications explicitly validate ES models against independent datasets (Chaplin-Kramer *et al.* 2019),  
50 although there are notable exceptions (Redhead *et al.* 2016; Sharps *et al.* 2017; Willcock *et al.* 2019). This  
51 is a particular issue as the results of location-specific validation (*e.g.* that performed during model  
52 development) may not be transferable to new locations (Redhead *et al.* 2016), or up-scalable to the regional  
53 and national extents over which ES model outputs are required to achieve the SDG (Willcock *et al.* 2016;  
54 Willcock *et al.* 2019). From a user and stakeholder perspective, not knowing the accuracy of the available  
55 ES models for the region of interest typically leads to either selection of a single suboptimal model – at  
56 worst leading to perverse decision-making – or a reluctance to use ES models altogether, causing an  
57 implementation gap between research, incorporation into policy and subsequent decision-making (Wong *et al.*  
58 *et al.* 2014; Willcock *et al.* 2016).  
59

60 Despite claims for predictive superiority of certain modelling techniques and platforms, independent  
61 evaluations have been unable to demonstrate the pre-eminence of any single approach. In fact, while more  
62 complex models on average perform better in terms of fit to validation data, the best-fit model varies  
63 regionally and often according to the validation data used (Sharps *et al.* 2017; Willcock *et al.* 2019; Willcock  
64 *et al.* 2020). So, if no single ES model is always the most accurate, how should a suitable approach be  
65 selected?  
66

67 Across the sciences, one solution to address uncertainty surrounding the accuracy of any single model is to  
68 use an ensemble of models (Araújo & New 2007; Willcock *et al.* 2020) – using individual models as  
69 replicates with different input parameters and boundary conditions (Araújo & New 2007; Dormann *et al.*  
70 2018). Variation among models in their assumptions and formats can result in large differences in  
71 predictions, in terms of predicted values and how they vary over space, especially when there is uncertainty  
72 as to the state and processes of the system being modelled (van Soesbergen & Mulligan 2018; Willcock *et al.*  
73 *et al.* 2019). Ensembles of models are hypothesised to have enhanced accuracy over individual models due to  
74 fewer overall errors in prediction by reducing the influence of idiosyncratic outcomes from single models  
75 (Araújo & New 2007; Dormann *et al.* 2018). Individual models rarely capture all potentially relevant  
76 processes or are often tuned to particular ecosystem characteristics. A combination of models might provide  
77 a more comprehensive coverage of processes and their forms, and avoids the chance of (unknowingly)  
78 selecting a model with a high prediction error at the location and scale of interest for a particular study  
79 (Willcock *et al.* 2020).  
80

81 Model ensembles are common in other disciplines – *e.g.* in niche modelling (Araújo & New 2007,  
82 Grenouillet *et al.* 2011), agroecology (Refsgaard *et al.* 2014), hydrology and water resources management  
83 (Wang *et al.* 2019; He *et al.* 2021), and climate and weather modelling (Knutti *et al.* 2013), as well as market  
84 forecasting (He *et al.* 2012). However, ensembles have been largely neglected in ES studies (Bryant *et al.*  
85 2018). The only current exception is the simplest ensemble approach (*i.e.* ‘committee averaging’ – taking  
86 the unweighted mean of a group of individual models per location –) which was applied to ES models in  
87 Sub-Saharan Africa, and gave higher accuracy in terms of fit to validation data (Willcock *et al.* 2020).  
88 Approaches that use more information might yield even more accurate estimates. Thus, here we explore the  
89 outstanding question of “what are the best ways to build ES model ensembles to realise the benefits such  
90 ensembles can bring to sustainability science?”

91  
92 Approaches to building model ensembles vary across disciplines, ranging from committee averaging  
93 (Marmion *et al.* 2009; Grenouillet *et al.* 2011) to complex Bayesian algorithms (Tebaldi & Knutti 2007).  
94 For example, species distribution models are generally deterministic statistical models; their fit to the data  
95 is often assessed with an accuracy metric and so ensembles are generally created using weighted averaging  
96 based on accuracy (Araújo & New 2007). By contrast, climate models are often treated as equal replicates  
97 with identical weights when making an ensemble (Tebaldi & Knutti 2007; Grenouillet *et al.* 2011) – we  
98 refer to such ensembles as ‘unweighted’. This difference may stem from the availability of suitable  
99 validation data, as well as different traditions. For example in species distribution models, biodiversity data  
100 are readily available and are used to train through cross-validation (Araújo & New 2007), whereas validation  
101 data on future climates obviously do not exist – although cross-validation against historic climate data is  
102 possible.

103  
104 As well as varying considerably in their underlying method, ES models often differ in the forms of their  
105 outputs, even when modelling the same ES (*e.g.* summed monetary value of the ES (de Groot *et al.* 2012)  
106 *vs.* specific biophysical predictions). By contrast, climate models generally have very similar forms of  
107 outputs. An important knowledge gap is therefore how to combine distinct ES model outputs as  
108 complementary inputs to provide a reliable ensemble. Outputs from different ES models can have different  
109 units and it is challenging to decide the relative weighting to place on each model. Models for a particular  
110 ES often have different structures, may include different processes, or may represent the same processes in  
111 different ways (Ochoa & Urbina-Cardona 2017). As a result, the different ES models will most likely not  
112 have equal accuracy, and so prediction errors (*i.e.* bias) may not be normally distributed among models  
113 (Dormann *et al.* 2018). If ES models had equal overall accuracies, unweighted averaging may provide a  
114 smoothing effect, reducing the impact of idiosyncratic outputs (*e.g.* at specific locations) of any particular  
115 model to reveal useful signals (Araújo & New 2007, Knutti *et al.* 2013; Diengdoh *et al.* 2020). In cases of  
116 varying overall accuracy, appropriate weighting of outputs based on model accuracy – *i.e.* models having  
117 unequal assigned weights – might re-adjust the distribution of prediction errors, and so improve the accuracy  
118 of the resulting ensemble (Refsgaard 2014; Dormann *et al.* 2018; Liu *et al.* 2020).

119  
120 However for ES, the lack of *a priori* validation data in many cases means that the distributions of accuracy  
121 among ES models are unknown. Furthermore, given that inferences about model accuracy at one location  
122 may not be transferable to others (Willcock *et al.* 2019), weighting using validation results from a separate  
123 study may not improve outcomes. Therefore where validation data are not available, the consensus among  
124 models could be used to weight their individual contribution to the ensemble value (Marmion *et al.* 2009;  
125 Grenouillet *et al.* 2011). This approach follows the logic that models whose output values are more different  
126 to those of the other models (*i.e.* are more distinct) are more likely to be incorrect. Therefore, weighting by  
127 consensus reduces the impact of outputs from more idiosyncratic models (*i.e.* those with extreme values,  
128 outliers or badly comparable processes) by comparison with the other models (Araújo & New 2007;  
129 Dormann *et al.* 2018), but does not exclude their information fully. The opposite may also be true – *i.e.*  
130 more distinct models are more accurate – for example in cases where more similar models have common  
131 inaccuracies.

132

133 Here, we implement 10 alternative ensemble methods, restricting ourselves to methods feasible for a wide  
 134 range of users, to evaluate whether weighting provides higher accuracy and if so which type of method  
 135 produces the most accurate predictions against validation data. We focus on two services, water supply and  
 136 carbon storage, in the United Kingdom. To support decision-making, we map the results for potential further  
 137 use, which are available via <https://doi.org/10.5285/a9ae773d-b742-4d42-ae42-2b594bae5d38>. We use  
 138 post-processing – specifically normalisation and per area correction – developed in earlier work (Willcock  
 139 *et al.* 2019; Willcock *et al.* 2020) to make outputs among models comparable.

## 141 2. Methods

142 We developed and validated unweighted average and weighted average ensembles of models for a  
 143 provisioning service (water supply; subsequently referred to as ‘water’) and a regulating service  
 144 (aboveground carbon storage; subsequently referred to as ‘carbon’), for which there is both a variety of  
 145 models available (Bagstad *et al.* 2013; Ochoa & Urbina-Cardona 2017; Willcock *et al.* 2019) and the  
 146 presence of accessible validation data. We applied the models and ensemble methods in the United Kingdom  
 147 (UK), for which there is a large quantity of reliable validation data; allowing us to assess ensemble  
 148 accuracies. We compared accuracy (*i.e.* fit to validation data) of these individual models with those of the  
 149 ensembles generated from them via multiple approaches, assessed if weighted ensembles were an  
 150 improvement on the unweighted mean-averaged ensemble, and identified the methods of weighting  
 151 ensembles that gave the highest accuracy.

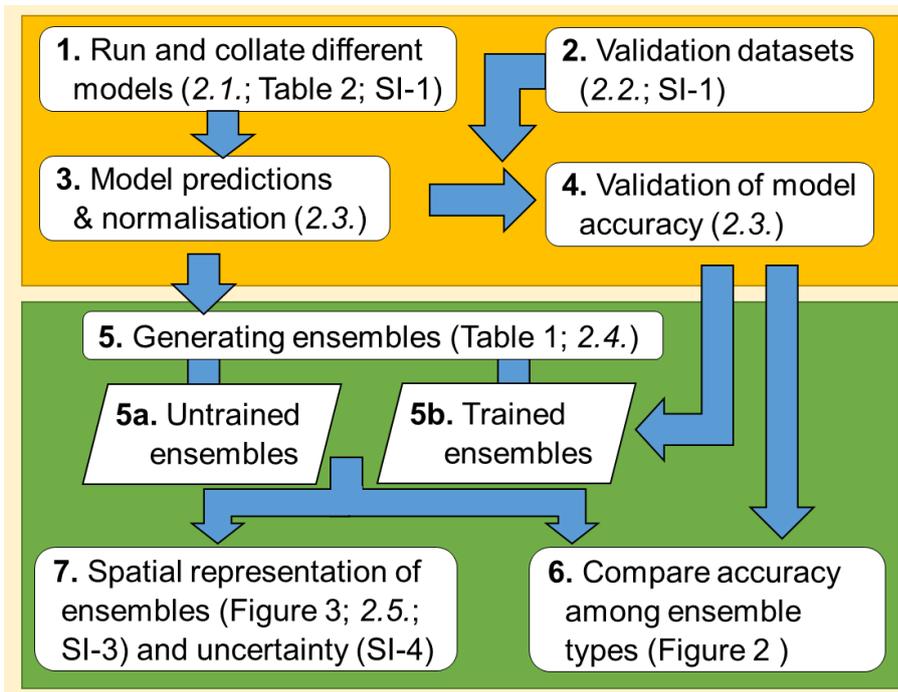
152  
 153 We modelled each ES at a 1 ha (100 × 100 m) resolution, and subsequently assessed performance of the  
 154 different ensemble approaches using weighting approaches we organised into three categories (Table 1):  
 155 deterministic consensus (*i.e.* always providing the same result), iterated consensus (*i.e.* using structured  
 156 trial-and-error approaches) and attribute-based (*e.g.* spatial resolution or distinctiveness). Finally, we  
 157 assessed the transferability of our UK results using independent data and models from a very different study  
 158 area – Sub-Saharan Africa (Willcock *et al.* 2019). We depict our overall process in Figure 1 in 7-steps. Our  
 159 calculations were performed using Matlab v7.14.0.739 and ArcMap 10.7.1, employing ArcPy coding for  
 160 loops. Relevant codes can be found at [github.com/EnsemblesTypes](https://github.com/EnsemblesTypes), with flow among codes explained in  
 161 SI-1-3.

162  
 163 **Table 1. Approaches used to calculate accuracy (A) and ensembles (B).** Ensemble approaches were  
 164 applied to the outputs of ten models for carbon storage and nine for water supply (see Table 2). For weighted  
 165 averaging, the procedure is described, and where applicable the Matlab tools used are mentioned; similar  
 166 regression tools are available in most statistical packages (further explanation is provided in SI-1). Trained  
 167 weighting (En-9 & En-10) uses validation data, whereas untrained weighting (En-3 to En-8) does not. En-1  
 168 and En-2 are unweighted average ensemble approaches, and En-3 to En-10 are weighted average  
 169 approaches; the latter comprising *deterministic* (En-3 & En-4), *iterated* (En-5, En-6 & En-10) and *attribute*  
 170 *weighted* (En-7 to En-9) techniques. With  $\omega_i$ : weight for model  $i$ ;  $E_{(x)}$ : the value of the ensemble;  $V_{(x)}$ : the  
 171 normalised validation value;  $Y_{i(x)}$  and  $Y_{j(x)}$ : the normalised value of model  $i$  or comparator  $j$  respectively, all  
 172 for selected spatial point  $x$ ; ( $y \neq x$ ) denoting a split dataset;  $C_{(i,j)}$ : the correlation coefficient between model  
 173  $i$  and  $j$ ; with  $n$  the # models,  $m$  the # spatial data points;  $n^g$ : the # models in distinctiveness group  $g$  (see SI-  
 174 1 for distinctiveness grouping).

Approach	Description	Details & Matlab Tool
<b>A. Accuracy approaches</b>		
• Spearman $\rho$	Correlation coefficient between ranked variables $V$ and $T$ .	$T$ is either $Y_i$ or $E$ , depending on ensemble method
• Inverse Deviance ( $D^\downarrow$ )	$D^\downarrow = 1 - \left(\frac{1}{m} \times \sum_x  X_{(x)} - T_{(x)} \right)$	$T_{(x)}$ is either $Y_{i(x)}$ or $\underline{E}_{(x)}$
<b>B. Ensemble approaches</b>		
<b>Unweighted Averaging:</b>		
En-1. Mean	$E_{(x)} = (\bar{Y}_i)_{(x)}$	

En-2. Median		$E_{(x)} = (\bar{Y}_i)_{(x)}$	Hypothesised to perform better than mean for skewed distributions.
<b>Untrained Weighted Ensembles: <math>E_{(x)} = \sum_i^n \left( \frac{\omega_i}{\sum_i^n \omega_i} \times Y_i \right)_{(x)}</math> with <math>\omega_i</math> following:</b>			
Deterministic consensus	En-3. PCA	$\omega_i$ = loadings of first Principal Component axis	Princomp-tool
	En-4. Correlation coefficients	$\omega_i = \frac{1}{n} \times \sum_j^n \frac{C_{(i,j)}}{\sqrt{C_{(i,i)} \times C_{(j,j)}}}$ , for all $j \in i$ with $C_{(i,j)} = \frac{1}{m-1} \times \sum_x^m \left( (Y_{i(x)} - \bar{Y}_i) \times (Y_{j(x)} - \bar{Y}_j) \right)$	
Iterated consensus	En-5. Regression to the median	$\bar{Y}_{(x)} \sim (\sum_i^n \omega_i Y_i)_{(x)}$	nlmefit-tool, maximising Log Likelihood
	En-6. Exhaustive leave-one-out cross-validation <sup>2</sup>	$Y_{j(x)} \sim \sum_{i \neq j}^n \omega_i Y_{i(x)}$ , for all $j \in i$ subsequently: $\omega_i = \frac{1}{n} \times \sum_i^n \left( \left( \frac{1}{n-1} \right) \times \sum_{i \neq j}^n \omega_{ij} \right)$	nlmefit-tool, maximising Log Likelihood
Attribute-based	En-7. Upweighted finer spatial resolution	$\omega_i = \frac{1}{\log_{10}(\text{spatial resolution})}$	Finer spatial resolution: smaller grid size in 1-dimensional meters (e.g. 25 m)
	En-8. Attribute weighting: distinctiveness	$\omega_i = \left( \frac{n^g}{n} \right)$ when upweighted with $n^g = i \in g$ $\omega_i = \left( \frac{n}{n^g} \right)$ when downweighted with $n^g = i \in g$	
<b>Trained Weighted Ensembles: <math>\omega</math>-transfer via jack-knife training</b>			
Attribute-based	En-9. Accuracy-weighted	$\omega_i = A_i$ , with $A_i(V_{(y \neq x)}, Y_{(y \neq x)})$	With $A$ , either Spearman $\rho$ or $D^{\downarrow}$ accuracy
Iterated consensus	En-10. Log-likelihood regressions	$V_{(y \neq x)} \sim (\sum_i^n \omega_i Y_i)_{(y \neq x)}$	Using nlmefit-tool, maximising Log Likelihood

176



177

178

179

180

181

182

183

184

**Figure 1.** Schematic representation of our ensemble analysis with arrows showing information flows. Numbers represent the steps with the method chapters indicated in italics, with respective detailing SIs; result figures are indicated. Parallelograms highlight the 10 ensemble approaches (Table 1), using models described in Table 2.

*2.1. Run and collate different models (step 1)*

185 We used outputs from 10 models for above ground carbon stocks based on per grid cell estimates, and  
186 outputs from nine models for annual water supply which provided accumulated flow estimates through  
187 specific pour points, either directly or through summation of run-off estimates per grid cell. We list these  
188 models in Table 2, including their output grid sizes (spatial resolution); we refer to SI-1-1 for full details,  
189 scales and supporting data. Acknowledging that model outputs have different units and sometimes model  
190 different constructs, we refer further to them in the general terms of carbon and water supply. Adhering to  
191 the aim of this paper, we do not compare individual model outputs, but focus on ensemble methods. All  
192 model outputs were set to the British National Grid transverse Mercator projection (EPSG 27700) with a  
193 0.9996 scale factor and units in metres. Not all models covered the whole of the UK, *e.g.* some excluded  
194 Northern Ireland or Scotland (see SI-1-1). Where applicable we corrected for this by using a standard error  
195 of means as  $(\frac{\sigma(x)}{\sqrt{n(x)}})$ , instead of standard deviation ( $\sigma$ ), with  $n$  the number of models per grid cell  $x$ . We  
196 collated models for this study according to their availability and to reflect different approaches to modelling  
197 ES.

198  
199  
200

**Table 2. Models and existing outputs used.** Full details, input data, post processing descriptions, and coverage are provided in SI-1-1. Model names are shown as acronyms and in full.

Model	Description	Grid size ( spatial resolution)	Model Type <sup>16</sup>
InVest v3.7.0 <sup>1†</sup> (Integrated Valuation of Ecosystem Services and Trade-offs)	Carbon module: above ground stocks	25 × 25 meters	Look-up table
	Water yield module: run-off per cell		Process
LPJ-GUESS <sup>2,3†</sup> (Lund-Potsdam-Jena General Ecosystem Simulator)	Vegetation biomass stocks per cell, mean for years 2009-2018	0.5° (≈ 46 × 46 km)	Process
	Water run-off per cell, mean for years 2009-2018		
LUCI <sup>4†</sup> (Land Utilisation Capability Indicator)	Above ground carbon stocks	10 × 10 meters	Look-up table
	Accumulated water run-off	5 × 5 meters	Process
\$-benefit transfer using The Economics of Ecosystems and Biodiversity database <sup>5,6†</sup>	Above ground carbon stock as monetary value	25 × 25 meters	Look-up table
	Water run-off as monetary value per cell		
Aqueduct v2.1 Total Blue Water <sup>7§</sup>	Accumulated water run-off	138 flow areas	Deterministic
ARIES k-Explorer <sup>8‡</sup> (Artificial Intelligence for Environment & Sustainability)	Joined above and below ground carbon stocks	1-hectare	Look-up table
Barredo <i>et al.</i> (2012) <sup>§</sup>	A European map of above ground biomass stocks	1 km <sup>2</sup>	Look-up table
Copernicus, Tree Cover Density <sup>9§</sup>	Proxy for carbon: tree Cover Density 2015 from MODIS satellite imagery.	20 × 20 meters	Deterministic
DECIPHeR <sup>10§</sup> (Dynamic fluxEs and ConnectIvity for Predictions of HydRology)	Accumulated water run-off through NRFA delineated catchment outlets, mean for years 1995-2015	387 catchments in common with validation	Process
Grid-to-Grid <sup>11§</sup>	Accumulated water run-off, mean for years 1995-2015	1 km <sup>2</sup>	Process
Henrys <i>et al.</i> (2016) <sup>§</sup>	Above ground carbon stocks	1 km <sup>2</sup>	Look-up table
Kindermann <i>et al.</i> (2008) <sup>§</sup>	A global map of above ground forest biomass stocks	1 hectare	Deterministic
National Forest Inventory (2018) <sup>12†</sup>	Woodland Land Cover Map <sup>15</sup> with above ground carbon stocks based on added Look-up table (Table. SI-1-4)	20 × 20 meters	Look-up table
Scholes Growth Days <sup>13,14†</sup>	Proxy for water run off per cell: # Days precipitation exceeds evapotranspiration	1 km <sup>2</sup>	Deterministic
WaterWorld v2 <sup>15‡</sup>	Accumulated water run-off	0.0083° (≈ 1 km <sup>2</sup> )	Process

201  
202

<sup>†</sup>Output generated for this work; <sup>‡</sup>online tool; <sup>§</sup>existing dataset; <sup>1</sup>Kareiva *et al.* (2011); <sup>2</sup>Smith *et al.* (2014); <sup>3</sup>Ahlström *et al.* (2015); <sup>4</sup>Thomas *et al.* (2020); <sup>5</sup>de Groot *et al.* (2012); <sup>6</sup>Costanza *et al.* (2014); <sup>7</sup>Gassert *et al.* (2015) <sup>8</sup>Martínez-López *et al.* (2019); <sup>9</sup>[land.copernicus.eu/tree-cover-density/status-maps/2015](http://land.copernicus.eu/tree-cover-density/status-maps/2015); <sup>10</sup>Coxon *et al.* (2019a; 2019b);

203 <sup>11</sup>Bell *et al.* (2018a; 2018b); <sup>12</sup>Forestry Commission (2018); <sup>13</sup>Scholes (1998); <sup>14</sup>Willcock *et al.* (2019); <sup>15</sup>Mulligan (2013); <sup>16</sup>following Ding & Bullock (2018), Willcock *et al.*  
204 (2019).  
205

## 2.2. Validation datasets (step 2)

Our carbon stock validation dataset was provided by Forest Research and comprises species inventories in all forest estates in England and Scotland in 2019 ([data-forestry.opendata.arcgis.com/](https://data-forestry.opendata.arcgis.com/); density shown in Figure 3; locations in Figure SI-1-2). In 201,143 forest compartments of varying size (mean: 4.4 hectares, median 1.6 hectares,  $\pm 22.1$ ), tree species, stand age and thinning regime were recorded for three vegetation layers. For each compartment and layer therein, the unique combination of stand age, thinning regime and tree species of the inventory data was searched in the UK Carbon Code tables ([woodlandcarboncod.org.uk](https://woodlandcarboncod.org.uk)) and life-time accumulated biomass was converted to total standing carbon per hectare estimates per compartment, with the layers summed per compartment (SI-1-2). Subsequently, compartments were spatially joined into 2078 polygons of ‘forest’ that were separated if more than 25 meters distance from each other.

Our water supply validation dataset comprised 519 hydrometric gauging stations from the National River Flow Archive of the UK (NRFA; [nrfa.ceh.ac.uk](https://nrfa.ceh.ac.uk)), with associated catchments representing a variety of sizes distributed across the whole of the UK (Figure 3). From the 1598 potential catchments in NRFA, we selected those that were  $>100 \text{ km}^2$  to get a robust mean run-off from the catchments. In cases where multiple gauging stations were found along the same river, based on name, only the largest was chosen to avoid pseudoreplication. An additional set of 41 Welsh catchments was included which did not meet this size criterion. Wales contains mainly small catchments due its geography – mountain ranges close to the sea – and so we selected catchments  $>25 \text{ km}^2$  to avoid this part of the UK being underrepresented. The data were polygons encompassing these catchments. Details are provided in SI-1-2.

## 2.3. Model predictions, normalisation (step 3) and validation of model accuracy (step 4)

For each individual model, predictions were obtained for each polygon in the validation dataset using the ArcGIS spatial analyst Zonal tool with a forced 2.5 m grid size environmental setting to minimise edge effects; *i.e.* all predicted values were obtained by resampling into  $2.5 \times 2.5 \text{ m}$  grid cells. In most cases the modelled value per polygon was obtained by taking the sum of all constituent grid cell values, corrected for both actual grid size and the resampling to 2.5 m. In the case of accumulated flow models, we corrected for potential small scale differences in flow routing among these models by taking the maximum flow value within both a 2 km range of the NRFA reported location of the gauging station and the polygon associated with that gauging station.

To ensure comparability among model outputs, we standardised by normalising among the outputs for each individual model and for the validation data-sets. Prior to this step all outputs were area corrected as either mean carbon stock – or proxy thereof – per hectare or water supply per hectare of catchment (with accumulated run-off estimates post-processed to give net run-off per cell; SI-1-1). This normalisation followed Willcock *et al.* (2019), and allowed us to address differences in units among models (such as monetary benefit transfer vs. satellite-based tree cover densities or run-off, and equalised carbon and biomass). To avoid impacts of extreme values without eliminating such data-points, we employed a double-sided Winsorising protocol for normalisation (Willcock *et al.* 2019; Verhagen *et al.* 2017), using the values associated to the 2.5% and 97.5% percentiles of number of datapoints to define the 0 and 1 values (values below or above these percentiles became 0 or 1 respectively). This winsorising normalisation protocol assumes outlier data are valid, but skewed values, in our case mainly by per area averaging, and corrects for this by compressing the variance tails rather than trimming them (Keselman *et al.* 2008; Erceg & Miroseovich 2008). Hence, we trade-off an even data distribution over the full 0-1 normalised range against the chance of having a true far outlier maximum (see SI-5 for a full investigation into the impact of the Winsorising protocol over standard normalisation for the validation data distribution). For each model, normalisation was done prior to creating ensembles.

For validation, we employed two accuracy measures (Willcock *et al.* 2019; Willcock *et al.* 2020), which are related to different aims in modelling ES (Table 1):

- 257 1) Comparing the rank order of predicted and validation data using Spearman  $\rho$ . This is relevant where  
 258 modelling is used to discover, for example, the most important locations for delivering an ES, or  
 259 conversely, those areas whose development may have least impact on ES delivery.  
 260 2) Ascertaining the absolute difference of each modelled value from its validation value using the inverse  
 261 of the deviance ( $D^\downarrow$ ). This is relevant where modelled values are important, *e.g.* when testing where ES  
 262 levels exceed a minimum threshold. We used the inverse of the deviance so that, like  $\rho$ , a higher value  
 263 indicated greater accuracy.

264  
 265 *2.4. Generate ensembles (step 5) and compare accuracy among ensemble types (step 6)*  
 266 We tested whether model ensembles were more accurate than the individual constituent models and which  
 267 approaches for creating ensembles were the most accurate in terms of fit to validation data. We created  
 268 ensembles using a range of methods, from the simplest calculation of an average value of the models at each  
 269 location ('unweighted averaged ensembles', *e.g.* Marmion *et al.* 2009, Grenouillet *et al.* 2011) to ensembles  
 270 with the contributions from different models weighted unequally ('weighted ensembles'), following  
 271 Dormann *et al.* (2018) (Table 1; further explanation and a model flow are provided in SI-1-3). We used  
 272 relatively straightforward approaches that would be feasible for a wide community of scientists and  
 273 decision-makers, and avoided more complex mathematical and/or statistical techniques such as Bayesian  
 274 networks (Bryant *et al.* 2018), which would require detailed specialist knowledge. Weights over all models  
 275 were normalised to sum to 1. Together with normalisation of the ensemble outputs (see above), this assured  
 276 equal scaling among all models and ensembles.

277  
 278 For unweighted average ensembles, we calculated both the mean and the median of modelled values at each  
 279 location as alternative measures of the central tendency which are differently affected by skew in the data  
 280 (Table 1, En-1 & En-2).

281  
 282 For weighted ensembles we calculated:  
 283 
$$E_{(x)} = \sum_i^n \left( \frac{\omega_i}{\sum_i^n \omega_i} \times Y_i \right)_{(x)}$$
 with positive weights  $\omega_i$  for model  $i$  of validation polygon  $x$ , weights  $\omega_i$  are  
 284 normalised to sum to 1,  $Y$  the modelled values for  $i$  per polygon (step 3), and  $n$  the total number  
 285 of models per service.  
 286

287 To determine  $\omega_i$ , the weighting value for each model  $i$ , we employed a range of methods that can be broadly  
 288 categorised as two main types of ensemble approach (untrained and trained), with further subdivision as:  
 289 deterministic consensus, iterated consensus, and attribute-based. The ensembles are listed as equations in  
 290 Table 1 (see SI-1-3 for further details).

- 291 1) Untrained ensembles (En-3 to En-8) represent a situation in which there is no validation data. To generate  
 292 uncertainty estimates allowing statistical comparison with the models and among ensembles we jack-  
 293 knifed (Araújo & New 2007; Refsgaard *et al.* 2014) with 50% of the spatial data polygons for 250 runs,  
 294 *i.e.* every run contained a new selection of half the dataset. We tested three approaches to produce the  
 295 ensembles:  
 296 - *Deterministic consensus* among models can be calculated using several approaches, including the fit  
 297 to a common consensus axis such as from a Principal Components Analysis (Marmion *et al.* 2009;  
 298 Grenouillet *et al.* 2011) or weighting by correlation coefficients (En-3 & En-4; ensemble numbering  
 299 follows Table 1).  
 300 - *Iterative approaches* might more accurately quantify consensus among models through using  
 301 structured trial-and-error (Dormann *et al.* 2018; Tebaldi & Knutti 2007). We use two regression  
 302 techniques: between the individual models and the median (En-5) and leave-one-out cross-validation  
 303 (En-6) following the suggestion in Dormann *et al.* (2018).  
 304 - One might *a priori* place value on a particular model attribute and use this to create weights (Englund  
 305 *et al.* 2017; Willcock *et al.* 2019; Brun *et al.* 2020; En-7, En-8 & En-9). For example, one could up-  
 306 or down-weight more distinct model types through a binary matrix of differences (En-8 & En-9; S1-

307 1-4) in land cover map used, grid-size, measured or modelled climate, model extent, presence of  
308 time-series, time step-size and model type (*i.e.* look-up table, deterministic or process based).  
309 Alternatively models that run at coarser spatial resolutions are penalised (En-7): smaller grid sizes  
310 are deemed more useful for decision-making (Willcock *et al.* 2016).

311 2) Trained ensembles (En-9 & En-10), as often used for species distribution models (*e.g.* Refsgaard *et al.*  
312 2014; Elith *et al.* 2011), represent a situation in which validation data are available from a similar region  
313 or part of the study area and so cannot be used to directly validate or substitute for the models in the  
314 study area, but can be used to weight these models. Here,  $\omega_i$  was trained with the validation data on a  
315 jack-knifed 50% of the dataset to achieve maximum accuracy (En-10) and subsequently  $\omega_i$  was  
316 transferred to the other half of the dataset. We used 250 such jack-knife runs (see above), with the same  
317 selections as above. Moreover, we included weighting by individual model accuracy (Marmion *et al.*  
318 2009; Liu *et al.* 2020) using the same jack-knife approach (En-9).

319  
320 After creating the ensembles, their accuracy was assessed following step 4 using the two measures (see 2.3):  
321 Spearman  $\rho$  and the inverse of the deviance ( $D^\dagger$ ). We assessed any improvement over the unweighted mean-  
322 averaged ensemble as the reference with pairwise t-tests against the null hypothesis of equal accuracy  
323 (Matlab *ttest-tool*). A similar analysis against the median-averaged ensemble as reference can be found in  
324 SI-2. To avoid spurious findings of significance through having a large number of replicates, we assessed  
325 improvement using bootstrapped tranches of 50 runs each with 250 replicates, and averaging the P-values.  
326 Since we used the same statistical test 12-times per service per accuracy estimate, we employed a full  
327 conservative Bonferroni correction; ( $\alpha = 0.05/12$ ) on the resulting average P-values. To compare the  
328 ensembles with the individual models we calculated per replicate the mean difference in accuracy among

329 all models ( $A_i$ ) against accuracy of an ensemble ( $A_E$ ) following:  $\left( \left( \sum_i^n \left( \frac{A_E}{A_i} - 1 \right) \right) \times \frac{1}{n} \right)$ , with n the number  
330 models and  $i$  an individual model.

331  
332 Steps 5 and 6 were repeated using independent data and models from a different study area (sub-Saharan  
333 Africa; Willcock *et al.* 2019) to investigate the transferability of the results presented here (Figure SI-2-2).

334

### 335 2.5. Spatial representation of ensembles and uncertainty (step 7)

336 To better support decision-making, we mapped our ES ensembles for the UK. For all the water ensembles,  
337 the mean normalised value across jack-knifed ensemble predictions per ensemble method were mapped as  
338 catchment polygons (step 5,  $N = 519$ ). For all carbon ensembles we mapped as 1 km<sup>2</sup> grid cells. Here, for  
339 each ensemble approach, the estimated weights as calculated for the validation polygons – mean averaged  
340 among jack-knife runs – were transferred to the full area, with the result aggregated to a 1 km<sup>2</sup> resolution  
341 based on the mean value among 1 hectare grid cells. In total, this carbon dataset has 253,802 cells that  
342 (partially) contain non-sea land cover. We transferred the weights calculated for the forests since running  
343 cross-validation approaches on over 250K data points would extremely time consuming to compute.  
344 However, since our validation data are only from forests/woodlands, we are aware of introducing a potential  
345 bias that could skew non-forested areas to lower values. Furthermore, we generated UK-scale maps of  
346 spatial variation in the differences among the untrained ensemble approaches, by calculating the standard  
347 error of the mean (SEM) among these spatial outputs. These maps are freely available online  
348 (<https://doi.org/10.5285/a9ae773d-b742-4d42-ae42-2b594bae5d38>), and spatial patterns of uncertainty are  
349 discussed in SI-4.

350

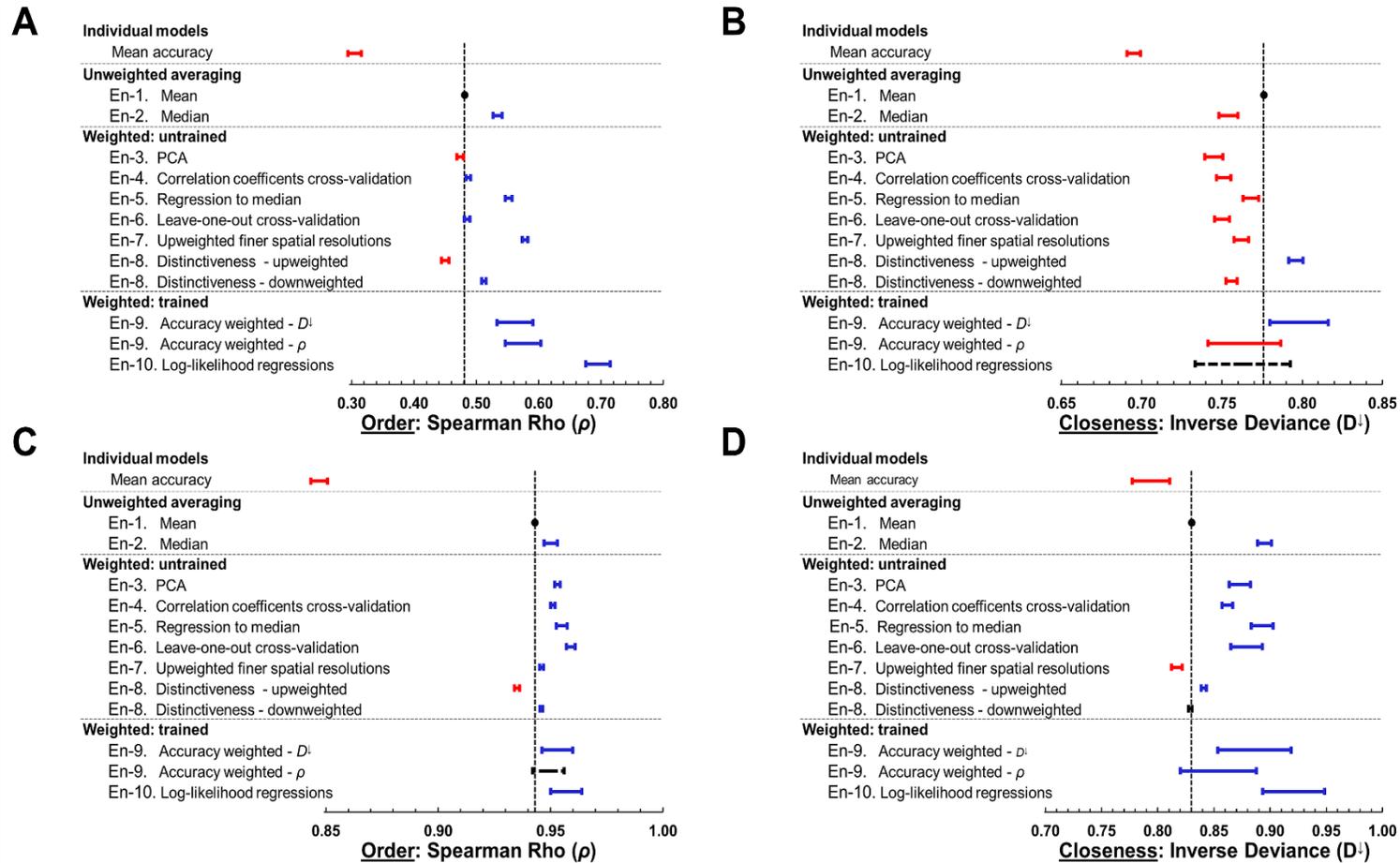
## 351 3. Results

352

### 353 3.1. Ensembles are more accurate than individual models

354 The average accuracy of individual models, represented by the mean of accuracy values taken across all  
355 models, was lower than that for any of the ensembles we created. The accuracy of the unweighted averaged  
356 ensembles (of modelled values at each location, *e.g.* ‘mean ensemble’) was appreciably higher than the

357 mean value for accuracy of the individual models for both carbon and water: 19%  $\pm$ 1.1% [sd] for  $\rho$  and  
358 12.1%  $\pm$ 0.5% for  $D^\downarrow$  improvement in fit to the validation data for carbon and 5.7%  $\pm$ 0.4% for  $\rho$  and 9.5%  
359  $\pm$ 1.7% for  $D^\downarrow$  for water (Figure 2). Untrained weighted ensembles showed large improvements – for most,  
360 larger than the unweighted ensembles – over the mean accuracy of the individual models of 17% to 27%  
361 ( $\rho$ ) and 7.6% to 15% ( $D^\downarrow$ ) for carbon (Figure 2A and B), and 5.3% to 6.5% ( $\rho$ ) and 7.7% to 18% ( $D^\downarrow$ ) for  
362 water (Figure 2C and D). In all cases, pairwise t-tests indicated highly significant differences between each  
363 ensemble and the mean value of accuracy of individual models (all  $P < 1E^{-10}$ ). Thus, creating an ensemble  
364 improves prediction accuracy against a randomly chosen individual model irrespective of the ensemble  
365 approach chosen.



366  
367  
368  
369  
370  
371  
372  
373  
374

**Figure 2. Accuracy of above ground carbon stock ensembles (10 models; A and B), and of water supply ensembles (9 models; C and D) against validation data.** The mean of accuracy values across the containing models – *i.e.* a randomly chosen model– is provided for comparison. For detail on the different ensemble types see Table 1 and SI-1-3. We show the average accuracy of 250 bootstrap runs with 50% of the dataset. The vertical dashed line indicates the reference accuracy of the unweighted mean-averaged ensemble (black dot, ‘mean ensemble’). Error bars indicate the standard deviation among runs in terms of proportional difference to the mean ensemble, calculated per bootstrap run as the difference in accuracy to the mean ensemble divided by the accuracy of the mean ensemble. The coefficient of variation among bootstraps for the mean carbon ensemble was 4% and 1%, for  $\rho$  and  $D^\downarrow$  respectively, and 1 % and 2% for water (not shown). **Blue** coloured ensemble accuracies are significantly higher than the unweighted mean ensemble (Bonferroni corrected  $\alpha = (0.05/12)$ ); **Red** coloured bars are significantly lower; **Black** dashed bars are not significantly different to the mean ensemble.

375 3.2. *Weighted ensembles are more accurate than unweighted ensembles*

376 All weighted ensembles, whether trained or untrained, significantly outperformed the reference unweighted  
377 mean ensemble (Figure 2), with the exception of  $D^\downarrow$  for carbon. In all cases, pairwise t-tests indicated these  
378 differences were highly significant ( $P < 1E^{-10}$ ; see Figure SI-2-1 for similar analyses against the median-  
379 averaged ensemble).

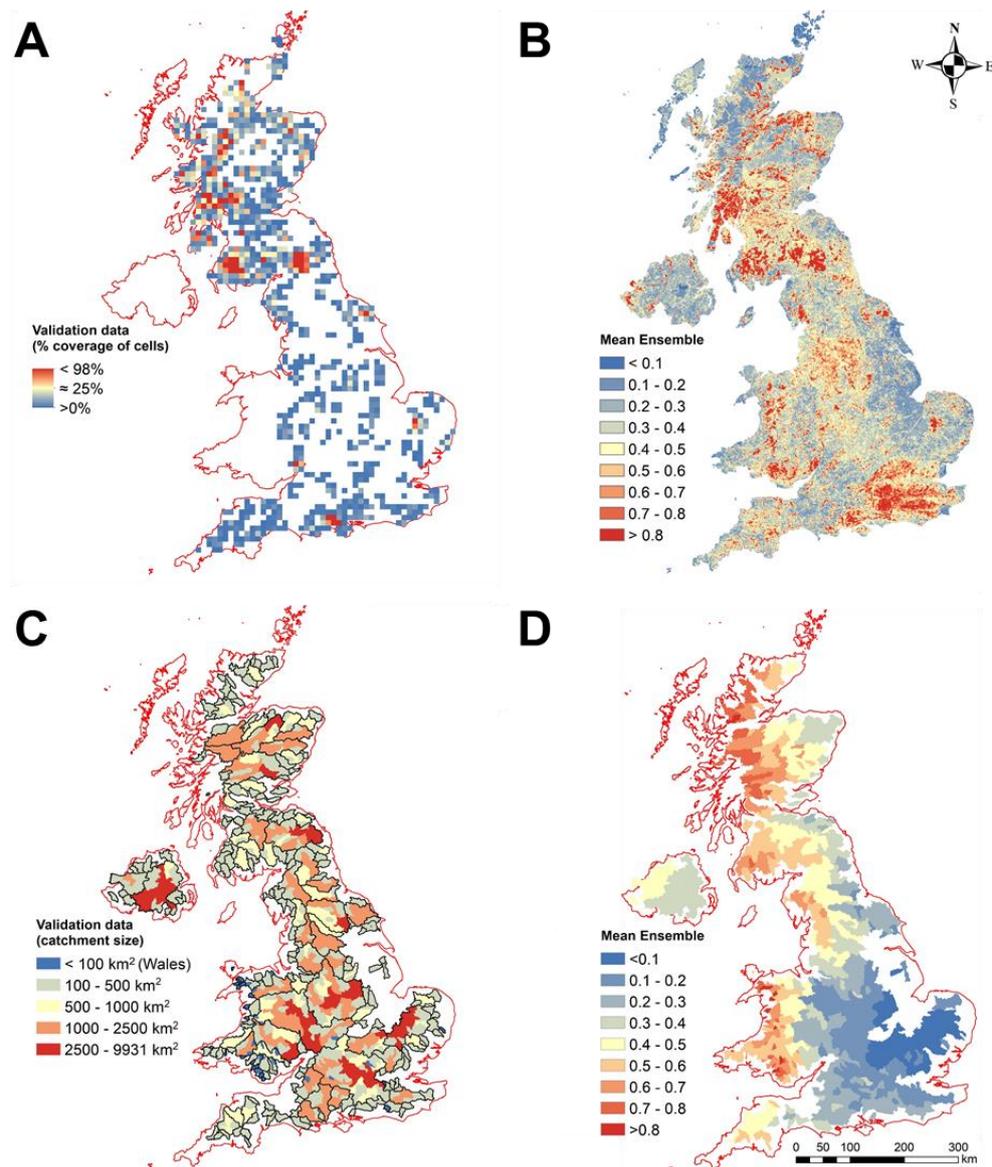
380  
381 For untrained weighted ensembles, prediction accuracy was elevated by up to  $4.8\% \pm 0.6\%$  for carbon  $\rho$   
382 (best: regression to median; Figure 2), with no improvement for carbon  $D^\downarrow$ , and  $0.8\% \pm 0.3\%$  and  $7.5\%$   
383  $\pm 1.1\%$  for water supply  $\rho$  and  $D^\downarrow$  respectively (regression to median; Figure 2). Conclusions as to the best  
384 model attributes to use for untrained weighting were dependent on the accuracy metric used ( $\rho$  or  $D^\downarrow$ ). By  
385 comparison to the unweighted mean ensembles, upweighting model outputs with finer spatial resolution  
386 improved  $\rho$  by up to  $6.6\% \pm 0.5\%$  and  $0.2\% \pm 0.1\%$  for carbon and water respectively but contrastingly  
387 decreased  $D^\downarrow$ . Upweighting more distinctive models was positive for  $D^\downarrow$  with  $2.5\% \pm 0.4\%$  and  $1.3\% \pm 0.3\%$   
388 greater accuracy compared to the unweighted mean ensemble for carbon and water supply respectively, but  
389 was negative for  $\rho$ . In summary, creating untrained weighted ensembles through iterative approaches was  
390 overall the most robust – particularly regression to the median (Table 1: En-5), showing greater accuracy  
391 than the unweighted mean-averaged ensembles in 3 out of 4 of our tests, and lower accuracy in 1 (Figure  
392 2).

393  
394 For trained weighting ensembles, using an iterative log-likelihood regression approach (Table 1: En-10) to  
395 establish weights elevated prediction accuracy compared to the unweighted mean ensemble by up to  $14.5\%$   
396  $\pm 2.6\%$  for carbon  $\rho$  (no improvement for carbon  $D^\downarrow$ ) and  $0.8\% \pm 0.7\%$  and  $11.1\% \pm 3.4\%$  for water supply  $\rho$   
397 and  $D^\downarrow$  respectively (Figure 2). Compared to such regressions, upweighting models with higher accuracy in  
398 the training set (accuracy-weighted ensembles; En-9; Figure 2) gave less improvement over the unweighted  
399 mean ensemble. Iteratively creating trained weighted ensembles using a log-likelihood regression approach  
400 (Table 1: En-10) was most robust – showing greater accuracy than the unweighted mean-averaged  
401 ensembles in 3 out of 4 of our tests, and is no worse in 1 (Figure 2).

402  
403 The reference unweighted mean ensembles for carbon and water are mapped for the UK in Figure 3. Maps  
404 for all other ensembles can be found in SI-3 and uncertainty among models and ensembles in SI-4. In  
405 accordance with *a priori* predictions, the uncertainty associated with selecting a single model was several  
406 times greater than that associated with selecting any single ensemble method for both ES. For carbon, the  
407 standard error of the means (SEM) among individual models per 1 km<sup>2</sup> grid cell (SEM =  $9.0\% \pm 2.8\%$ , SI-  
408 4) was ca. 3.5-times larger than among ensembles (SEM =  $2.5\% \pm 1.1\%$ ). Similarly, the SEM among  
409 individual water models per watershed (SEM =  $7.8\% \pm 3.4\%$ , SI-4) was substantially greater than among  
410 ensembles (SEM =  $1.3\% \pm 0.7\%$ ). In SI-4 we investigate spatial drivers for this uncertainty, discussing these  
411 patterns at length.

412  
413 We validated the robustness of our results using independent data and models from a different area (Sub-  
414 Saharan Africa; Willcock *et al.* 2019), which gave similar results of weighted ensembles outperforming the  
415 reference mean ensemble (Figure SI-2-2).

416



417  
 418 **Figure 3. Spatial distribution of validation points and the reference mean ecosystem service value. A**  
 419 **the Distribution of 2078 carbon validation forests as coverage of 10 × 10 km cells – many individual forest**  
 420 **fragments would be too small to be clear at this scale, see SI SI-1-2 –, white cells are empty. B** the reference  
 421 **unweighted mean ensemble of carbon across 10 models, normalised on scale 0-1. C** the 519 catchments  
 422 **used for water validation and ensemble calculations coloured by their size – smaller watersheds that overlap**  
 423 **larger ones are displayed on top; lines show underlying largest catchment level. D** the reference unweighted  
 424 **mean ensemble of water supply across 9 models, normalised on scale 0-1. All maps here, in SI-3 (all**  
 425 **ensembles) and SI-4 (uncertainty) could support landscape decisions in the UK and are available via**  
 426 <https://doi.org/10.5285/a9ae773d-b742-4d42-ae42-2b594bae5d38>.

427

#### 428 4. Discussion

429 We have shown that predictions from ensembles of models have substantially higher accuracy than a  
 430 randomly selected single ES model, and especially that weighting approaches increase ensemble accuracy.  
 431 Finding increased performance through use of ensemble approaches is common in other fields. For example,  
 432 the increased accuracy of ensemble species distribution models ranges from 1-2% (Crossman *et al.* 2012;  
 433 Abrahms *et al.* 2019) to 12% (Grenouillet *et al.* 2011), although an increase is not universal (Hao *et al.*  
 434 2020). Similarly, 2% accuracy increases were found for market forecasting ensembles (He *et al.* 2012), and  
 435 neural network ensemble averaging resulted in up to 7% improvements in accuracy (Inoue & Narisha 2000).

436  
437 Specific to ES, unweighted averaged ensembles have been shown to be 5.0–6.1% more accurate than  
438 individual models (Willcock *et al.* 2020). Our improvements with ES ensembles are at minimum 5%-17%,  
439 suggesting substantial differences among models in their adequacy (Dormann *et al.* 2018), but also that  
440 ensemble approaches that use more information offer greater increases in accuracy. We found that taking  
441 the median generally outperforms a mean ensemble, probably because the latter is more influenced by  
442 outliers. Our results provide evidence that weighted ES ensembles created using consensus techniques  
443 produce more accurate outputs than unweighted ensembles. This finding is supported by our additional  
444 analysis using independent models and data from Sub-Saharan Africa (in a biome with very different  
445 climatic and soil characteristics; SI-2), suggesting our findings may be generalisable, although investigating  
446 this specifically (e.g., for different ES, regions and validation datasets) is an important avenue for future  
447 research.

448  
449 Predictions from models, including those from ES models, are all potentially biased in direction and amount  
450 because of their underlying assumptions. These biases could differ among models due to their specific  
451 construction. Therefore, models are likely to differ in their accuracy when compared to reality (Dormann *et al.*  
452 *et al.* 2018). The improvement in accuracy when using ensembles, as we have shown here, is referred to as a  
453 ‘portfolio effect’ by which a (weighted) combination of replications of possible states of a system suppresses  
454 idiosyncratic differences and provides a more reliable average estimate (Thibaut & Connolly 2013;  
455 Dormann *et al.* 2018; Lewis *et al.* 2021). However, this effect is lessened if models share similar  
456 assumptions and, therefore, concomitant biases – highlighting the importance of including multiple model  
457 outputs (Ding & Bullock 2018) and, where data are available, model validation (Willcock *et al.* 2019). In  
458 particular, the use of models not usually packaged as ES models – such as LPJ-GUESS – might help with  
459 increasing the variety of inputs for ensembles. If some models systematically overestimate and other models  
460 underestimate, averaging delivers smaller prediction errors when models are weighted (Dormann *et al.*  
461 2018). Hence, the resulting weighted ensemble is more accurate than most individual models and  
462 unweighted approaches (Marmion *et al.* 2009, Grenouillet *et al.* 2011); see Dormann *et al.* (2018) for  
463 theoretical explorations.

464  
465 We have shown the general potential of weighting to re-balance the contribution of different ES models,  
466 but also find that some weighting approaches seem more suitable. Specifically, structured trial-and-error  
467 iterative approaches may more accurately maximise consensus among models than deterministic approaches  
468 (Dormann *et al.* 2018; Gobeyn *et al.* 2019). The PCA and correlation coefficient approaches (Table 1: En-  
469 3 & En-4) deterministically assess consensus among individual models. By contrast, regression to the  
470 median, leave-one-out cross validation, and log-likelihood approaches (Table 1: En-5, En-6, En-10) are  
471 examples of iterative processes that optimise for the highest level of consensus in full parameter space  
472 (Dormann *et al.* 2018). Attribute-based approaches as used by Masson & Knutti (2011) and Willcock *et al.*  
473 (2019) (e.g. weighting by model distinctiveness or grid size; Table 1: En-7 and En-8) produce conflicting  
474 results. Model attributes such as these may not correctly describe why model outputs vary, or capture their  
475 complexity (Willcock *et al.* 2019; Brun *et al.* 2020) and so weighting by among-model agreement produces  
476 more accurate ensemble outputs. One might expect accuracy-weighted ensembles (Table 1: En-9) to  
477 perform best. However, model accuracy can be location specific and poorly transferable elsewhere – even  
478 with similar model accuracy, some grid cells may be well represented by some models and less by others  
479 (Graham *et al.* 2008; Marmion *et al.* 2009; Zulian *et al.* 2018). As a result accuracy-derived weights show  
480 high uncertainty in areas where training data were not available (i.e. non-forested areas; SI-4), likely because  
481 of over-fitting to areas with available data (i.e. forests/woodlands) producing correlative patterns that  
482 explain other areas less well. In SI-4, we investigated environmental and spatial drivers of uncertainty  
483 among predictions. Broadly, these supplementary results show that carbon models and ES ensembles are  
484 less accurate in urban areas. We also find that ensembles for water are less accurate in areas of high rainfall,  
485 seasonality and rugosity (see SI-4 for full details). That said, as uncertainty among ES ensembles is almost  
486 4-times lower than among individual models, this suggests less need to make the ‘right choice’ of method

487 when selecting an ensemble approach. Thus, although there is some chance of picking a superior individual  
488 model (Willcock *et al.* 2018), the risk of a sub-optimal prediction is substantially lowered by applying any  
489 ensemble method and this risk is further reduced when a weighted ensemble is used.

490  
491 Our results should serve as a ‘call to arms’ for ES researchers and practitioners to increasingly use ensembles  
492 of models to support decision-making for sustainability. Using an individual ES model is fraught with  
493 concerns as *a priori* it is not known which is the most accurate and choosing only one model can, at worst,  
494 result in perverse decisions (Willcock *et al.* 2019). Deriving decisions from an ensemble of ES models  
495 provides an improvement over using one model for any location (which may be large or small, depending  
496 on the local context and the models used), but also more consistency over space, as model accuracy varies  
497 spatially (see results in SI-4). Therefore, using ensemble approaches, and especially weighted ensembles,  
498 would increase credibility and so help reduce the implementation gap between research and policy- and  
499 decision-making (Wong *et al.* 2014; Willcock *et al.* 2016). We acknowledge the lack of standardised metrics  
500 across models and limited computational and financial resources that could restrict the uptake of ensembles  
501 – indeed, many practitioners only run a single model. However, given the errors associated with single  
502 models (this paper; Willcock *et al.* 2020; Eigenbrod *et al.* 2010), we argue that a single model is inadequate,  
503 although more complex models are sometimes more accurate (Willcock *et al.* 2019). The most complex (a  
504 priori best) ES models require substantial inputs (i.e. data, computational power, subscription fees, and staff  
505 time), and so running multiple models – whilst requiring additional resources – results in a large gain per  
506 extra unit resource. For example, as even untrained weighted ensembles developed using iterative  
507 approaches (e.g. regression to the median, leave-one-out cross validation) enable a 3-fold reduction in  
508 variation, such an ensemble approach seems a reasonable minimum standard for ES modelling – striking  
509 the right balance between feasibility and robustness (Willcock *et al.* 2016). Whilst such ensembles will be  
510 outperformed by the best-performing individual models, these cannot be identified without running multiple  
511 models – a ‘Catch-22’ (Willcock *et al.* 2019). Thus, we recommend that multiple models be developed for  
512 ES where they are lacking (e.g. cultural services; Martínez-Harms and Balvanera, 2012; Wong *et al.* 2014),  
513 and that those with access to sufficient resources to run multiple models ensure the ensemble outputs are  
514 freely available, making the use of these ensembles more feasible and accessible for all (Willcock *et al.*  
515 2020).

516

## 517 **5. Conclusion**

518 We show that in situations with no *a priori* validation evidence guiding model selection, predictions from  
519 ensembles of models have a higher accuracy than selecting an individual model by chance. Weighted  
520 averaging further improves accuracy, suppressing idiosyncratic differences through producing consensus  
521 (Araújo & New 2007; Dormann *et al.* 2018). Doing so not only elevates accuracy but substantially decreases  
522 uncertainty among ensemble approaches compared to uncertainty among models, a further indication of  
523 increased fit to reality (Chaplin-Kramer *et al.* 2019; Willcock *et al.* 2020). In summary, even if a less  
524 accurate ensemble weighting approach is used, one would on average have lower uncertainty than selecting  
525 an individual model by chance. Thus, particularly when validation data are not available, we recommend  
526 the use of weighted ensembles in ES research to substantially reduce uncertainty and to support robust  
527 decision-making for sustainable development.

528

## 529 **References**

- 530 Abrahms, B. *et al.* (2019). Dynamic ensemble models to predict distributions and anthropogenic risk  
531 exposure for highly mobile species. *Divers. Distrib.* **25**, 11821193.  
532 <https://doi.org/10.1111/ddi.12940>
- 533 Ahlström, A. *et al.* (2015). Carbon cycle. The dominant role of semi-arid ecosystems in the trend and  
534 variability of the land CO<sub>2</sub> sink. *Science* **348**, 895–899. <https://doi.org/10.1126/science.aaa1668>
- 535 Araújo, M.B. & New, M. (2007). Ensemble forecasting of species distributions. *Trends Ecol. Evol.* **22**,  
536 42–47. <https://doi.org/10.1016/j.tree.2006.09.010>.

537 Bagstad, K.J. *et al.* (2013). A comparative assessment of decision-support tools for ecosystem services  
538 quantification and valuation. *Ecosyst. Serv.* **5**, 27–39. <https://doi.org/10.1016/j.ecoser.2013.07.004>  
539 Barredo, J.I. *et al.* (2012). *A European map of living forest biomass and carbon stock.* (European  
540 Commission, Joint Research Centre). [https://op.europa.eu/en/publication-detail/-](https://op.europa.eu/en/publication-detail/-/publication/b9345574-a96f-4417-87ed-1a85d2252834/language-en)  
541 [/publication/b9345574-a96f-4417-87ed-1a85d2252834/language-en](https://op.europa.eu/en/publication-detail/-/publication/b9345574-a96f-4417-87ed-1a85d2252834/language-en)  
542 Bell, V.A. *et al.* (2018a). The MaRIUS- G2G datasets: Grid- to- Grid model estimates of flow and  
543 soil moisture for Great Britain using observed and climate model driving data. *Geosci. Data J.* **5**,  
544 63–72. <https://doi.org/10.1002/gdj3.55>  
545 Bell, V.A. *et al.* (2018b). *Grid-to-Grid model estimates of monthly mean flow and soil moisture for*  
546 *Great Britain (1891 to 2015): observed driving data [MaRIUS-G2G-Oudin-monthly].* [Data Set]  
547 (NERC Environmental Information Data Centre). [https://doi.org/10.5285/f52f012d-9f2e-42cc-](https://doi.org/10.5285/f52f012d-9f2e-42cc-b628-9cdea4fa3ba0)  
548 [b628-9cdea4fa3ba0](https://doi.org/10.5285/f52f012d-9f2e-42cc-b628-9cdea4fa3ba0)  
549 Brun, P. *et al.* (2020). Model complexity affects species distribution projections under climate  
550 change. *J. Biogeogr.* **47**, 130–142. <https://doi.org/10.1111/jbi.13734>  
551 Bryant, B.P. *et al.* (2018). Transparent and feasible uncertainty assessment adds value to applied  
552 ecosystem services modeling. *Ecosyst.Serv.* **33**, 103–109.  
553 <https://doi.org/10.1016/j.ecoser.2018.09.001>  
554 Chaplin-Kramer, R. *et al.* (2019). Global modeling of nature’s contributions to people. *Science* **366**,  
555 255–258. <https://science.sciencemag.org/content/366/6462/255.abstract>  
556 Costanza, R. *et al.* (2014). Changes in the global value of ecosystem services. *Glob. Environ.*  
557 *Change* **26**, 152–158. <https://doi.org/10.1016/j.gloenvcha.2014.04.002>  
558 Costanza, R. *et al.* (2017). Twenty years of ecosystem services: how far have we come and how far do  
559 we still need to go? *Ecosyst. Serv.* **28**, 1–16. <https://doi.org/10.1016/j.ecoser.2017.09.008>  
560 Coxon, G. *et al.* (2019a). DECIPHeR v1: Dynamic fluxEs and ConnectIvity for Predictions of  
561 HydRology. *Geosci. Model Dev.* **12**, 2285–2306. <https://doi.org/10.5194/gmd-12-2285-2019>  
562 Coxon, G. *et al.* (2019b). *DECIPHeR model estimates of daily flow for 1366 gauged catchments in*  
563 *Great Britain (1962-2015) using observed driving data.* [Data Set] (NERC Environmental  
564 Information Data Centre). <https://doi.org/10.5285/d770b12a-3824-4e40-8da1-930cf9470858>  
565 Crossman, N.D., Bryan, B.A. & Summers, D.M. (2012). Identifying priority areas for reducing species  
566 vulnerability to climate change. *Divers. Distrib.* **18**, 60–72. [https://doi.org/10.1111/j.1472-](https://doi.org/10.1111/j.1472-4642.2011.00851.x)  
567 [4642.2011.00851.x](https://doi.org/10.1111/j.1472-4642.2011.00851.x)  
568 Diengdoh, V.L. *et al.* (2020). A validated ensemble method for multinomial land-cover  
569 classification. *Ecol. Inform.* **56**, 101065. <https://doi.org/10.1016/j.ecoinf.2020.101065>  
570 Ding, H. & Bullock, J.M. (2018). *A Guide to Selecting Ecosystem Service Models for Decision-*  
571 *Making: Lessons from Sub-Saharan Africa.* (World Resources Institute). [wri.org/publication/guide-](http://wri.org/publication/guide-selecting-ecosystem-service)  
572 [selecting-ecosystem-service](http://wri.org/publication/guide-selecting-ecosystem-service)  
573 Dormann, C.F. *et al.* (2018). Model averaging in ecology: a review of Bayesian, information-theoretic,  
574 and tactical approaches for predictive inference. *Ecol. Monogr.* **88**, 485–504.  
575 <https://doi.org/10.1002/ecm.1309>  
576 Eigenbrod, F. *et al.* (2010) The impact of proxy- based methods on mapping the distribution of  
577 ecosystem services. *J. Appl. Ecol.* **47.2**, 377-385.  
578 Elith, J. *et al.* (2011). A statistical explanation of MaxEnt for ecologists. *Divers. Distrib.* **17**, 43–57.  
579 <https://doi.org/10.1111/j.1472-4642.2010.00725.x>  
580 Englund, O., Berndes, G. & Cederberg, C. (2017). How to analyse ecosystem services in landscapes—A  
581 systematic review. *Ecol. Indic.* **73**, 492–504. <https://doi.org/10.1016/j.ecolind.2016.10.009>  
582 Erceg-Hurn, D.M. & Mirosevich, V.M. (2008). Modern robust statistical methods: an easy way to  
583 maximize the accuracy and power of your research. *Am. Psychol.* **63**, 591–601.  
584 <http://dx.doi.org/10.1037/0003-066X.63.7.591>  
585 Forestry Commission, United Kingdom. (2018). *National Forest Inventory Woodland GB 2018.* [Data  
586 Set] (Forestry Commission Open Data). [http://data-](http://data-forestry.opendata.arcgis.com/datasets/d3d7bfba1cba4a3b83a948f33c5777c0_0)  
587 [forestry.opendata.arcgis.com/datasets/d3d7bfba1cba4a3b83a948f33c5777c0\\_0](http://data-forestry.opendata.arcgis.com/datasets/d3d7bfba1cba4a3b83a948f33c5777c0_0)

588 Gassert, F. *et al.* (2015). *Aqueduct Global Maps 2.1*. [Data Set] (World Resources Institute).  
589 <https://www.wri.org/resources/data-sets/aqueduct-global-maps-21-data>  
590 Gobeyn, S. *et al.* (2019). Evolutionary algorithms for species distribution modelling: A review in the  
591 context of machine learning. *Ecol. Modell.* **392**, 179–195.  
592 <https://doi.org/10.1016/j.ecolmodel.2018.11.013>  
593 Graham, C.H. *et al.* (2008). The influence of spatial errors in species occurrence data used in  
594 distribution models. *J Appl. Ecol.* **45**, 239–247. <https://doi.org/10.1111/j.1365-2664.2007.01408.x>  
595 Grenouillet, G. *et al.* (2011). Ensemble modelling of species distribution: the effects of geographical  
596 and environmental ranges. *Ecography* **34**, 9–17. <https://doi.org/10.1111/j.1600-0587.2010.06152.x>  
597 Griggs, D. *et al.* (2013). Sustainable development goals for people and planet. *Nature* **495**, 305–307.  
598 <https://doi.org/10.1038/495305a>  
599 de Groot, R. *et al.* (2012). Global estimates of the value of ecosystems and their services in monetary  
600 units. *Ecosyst. Serv.* **1**, 50–61. <https://doi.org/10.1016/j.ecoser.2012.07.005>  
601 Hao, T. *et al.* (2020). Testing whether ensemble modelling is advantageous for maximising predictive  
602 performance of species distribution models. *Ecography* **43**, 549–558.  
603 <https://doi.org/10.1111/ecog.04890>  
604 He, X. . *et al.* (2021). Climate-informed hydrologic modeling and policy typology to guide managed  
605 aquifer recharge. *Science Advances* **7**, p.eabe6025. <https://doi.org/10.1126/sciadv.abe6025>  
606 He, K., Yu, L. & Lai, K.K. (2012). Crude oil price analysis and forecasting using wavelet decomposed  
607 ensemble model. *Energy* **46**, 564–574. <https://doi.org/10.1016/j.energy.2012.07.055>  
608 Henrys, P.A., Keith, A. & Wood, C.M. (2016). *Model estimates of aboveground carbon for Great*  
609 *Britain*. [Data Set] (NERC Environmental Information Data  
610 Centre). <https://doi.org/10.5285/9be652e7-d5ce-44c1-a5fc-8349f76f5f5c>  
611 Inoue, H. & Narihisa, H. (2000) in *Knowledge Discovery and Data Mining. Current Issues and New*  
612 *Applications* (eds Terano, T, Liu, H. & Chen, A.L.P.) 177-180 (Springer).  
613 <https://link.springer.com/book/10.1007/3-540-45571-X>  
614 Kareiva, P. *et al.* (2011). *Natural Capital: Theory and Practice of Mapping Ecosystem Services*.  
615 (Oxford University Press).  
616 [https://oxford.universitypressscholarship.com/view/10.1093/acprof:oso/9780199588992.001.0001/a](https://oxford.universitypressscholarship.com/view/10.1093/acprof:oso/9780199588992.001.0001/acprof-9780199588992)  
617 [cprof-9780199588992](https://oxford.universitypressscholarship.com/view/10.1093/acprof:oso/9780199588992.001.0001/acprof-9780199588992)  
618 Keselman, H. J. *et al.* (2008). A generally robust approach for testing hypotheses and setting  
619 confidence intervals for effect sizes. *Psychol. Methods* **13**, 110–129.  
620 <https://doi.apa.org/doi/10.1037/1082-989X.13.2.110>  
621 Kindermann, G.E. *et al.* (2008). A global forest growing stock, biomass and carbon map based on FAO  
622 statistics. *Silva Fennica* **42**, 397–396. <http://pure.iiasa.ac.at/id/eprint/8616/>  
623 Knutti, R., Masson, D. & Gettelman, A. (2013). Climate model genealogy: Generation CMIP5 and  
624 how we got there. *Geophys. Res. Lett.* **40**, 1194–1199. <https://doi.org/10.1002/grl.50256>  
625 Lewis, K.A. . *et al.* (2021). Using multiple ecological models to inform environmental decision-  
626 making. *Front. Mar. Sci.* **8**, 283. <https://doi.org/10.3389/fmars.2021.625790>  
627 Liu, D., Li, T. & Liang, D. (2020). An integrated approach towards modeling ranked weights. *Comput.*  
628 *Ind. Eng.* **147**, 106629. <https://doi.org/10.1016/j.cie.2020.106629>  
629 Malinga, R. *et al.* (2015). Mapping ecosystem services across scales and continents—A review. *Ecosyst.*  
630 *Serv.* **13**, 57–63. <https://doi.org/10.1016/j.ecoser.2015.01.006>  
631 Marmion, M. *et al.* (2009). Evaluation of consensus methods in predictive species distribution  
632 modelling. *Divers. Distrib.* **15**, 59–69. <https://doi.org/10.1111/j.1472-4642.2008.00491.x>  
633 Martínez-Harms, M.J. & Balvanera, P. (2012). Methods for mapping ecosystem service supply: a  
634 review. *Int. J. Biodivers. Sci. Ecosyst. Serv. Manag.* **8**, 17-25.  
635 Martínez-López, J. *et al.* (2019). Towards globally customizable ecosystem service models. *Sci. Total*  
636 *Environ.* **650**, 2325–2336. <https://doi.org/10.1016/j.scitotenv.2018.09.371>  
637 Masson, D. & Knutti, R. (2011). Climate model genealogy. *Geophys. Res. Lett.* **38**. L08703.  
638 <https://doi.org/10.1029/2011GL046864>

639 Mulligan M. (2013). WaterWorld: a self-parameterising, physically based model for application in  
640 data-poor but problem-rich environments globally. *Hydrol. Res.* **44**, 748–69.  
641 <https://doi.org/10.2166/nh.2012.217>

642 Ochoa, V. & Urbina-Cardona, N. (2017). Tools for spatially modeling ecosystem services: Publication  
643 trends, conceptual reflections and future challenges. *Ecosyst.Serv.* **26**, 155–169.  
644 <https://doi.org/10.1016/j.ecoser.2017.06.011>

645 Pascual, U. *et al.* (2017). Valuing nature’s contributions to people: the IPBES approach. *Curr. Opin.*  
646 *Environ. Sustain.* **26–27**, 7–16. <https://doi.org/10.1016/j.cosust.2016.12.006>

647 Redhead, J.W. *et al.* (2016). Empirical validation of the InVEST water yield ecosystem service model  
648 at a national scale. *Sci. Total Environ.* **569**, 1418–1426 (2016).  
649 <https://doi.org/10.1016/j.scitotenv.2016.06.227>

650 Refsgaard, J.C. *et al.* (2014). A framework for testing the ability of models to project climate change  
651 and its impacts. *Clim. Change* **122**, 271–282. <https://doi.org/10.1007/s10584-013-0990-2>

652 Scholes, R.J. (1998). *The South African I: 250 000 maps of areas of homogeneous grazing potential.*  
653 (CSIR, South Africa). No internet reference

654 Sharps, K. *et al.* (2017). Comparing strengths and weaknesses of three ecosystem services modelling  
655 tools in a diverse UK river catchment. *Sci. Total Environ.* **584**, 118–130.  
656 <https://doi.org/10.1016/j.scitotenv.2016.12.160>

657 Smith, B. *et al.* (2014). Implications of incorporating N cycling and N limitations on primary  
658 production in an individual-based dynamic vegetation model. *Biogeosciences* **11**, 2027–2054.  
659 <https://d-nb.info/1121909426/34>

660 van Soesbergen, A. & Mulligan, M. (2018). Uncertainty in data for hydrological ecosystem services  
661 modelling: Potential implications for estimating services and beneficiaries for the CAZ Madagascar.  
662 *Ecosyst. Serv.* **33**, 175–186. <https://doi.org/10.1016/j.ecoser.2018.08.005>

663 Suich, H., Howe, C. & Mace, G. (2015). Ecosystem services and poverty alleviation: A review of the  
664 empirical links. *Ecosyst. Serv.* **12**, 137–147. <https://doi.org/10.1016/j.ecoser.2015.02.005>

665 Tebaldi, C. & Knutti, R. (2007). The use of the multi-model ensemble in probabilistic climate  
666 projections. *Philos. Trans. R. Soc. A Math. Phys. Eng. Sci.* **365**, 2053–2075.  
667 <https://doi.org/10.1098/rsta.2007.2076>

668 Thibaut, L.M. & Connolly, S.R. (2013). Understanding diversity–stability relationships: towards a  
669 unified model of portfolio effects. *Ecol. Lett.* **16**, 140–150. <https://doi.org/10.1111/ele.12019>

670 Thomas, A. *et al.* (2020). Fragmentation and thresholds in hydrological flow- based ecosystem  
671 services. *Ecol. Appl.* **30**, e02046. <https://doi.org/10.1002/eap.2046>

672 UKNEA. (2011). *The UK National Ecosystem Assessment: Synthesis of the Key Findings.* (UNEP-  
673 WCMC, Cambridge). [https://www.unep-wcmc.org/resources-and-data/UK-national-ecosystem-](https://www.unep-wcmc.org/resources-and-data/UK-national-ecosystem-assessment)  
674 [assessment](https://www.unep-wcmc.org/resources-and-data/UK-national-ecosystem-assessment)

675 Verhagen, W. *et al.* (2017). Use of demand for and spatial flow of ecosystem services to identify  
676 priority areas. *Conserv. Biol.* **31**, 860–871. <https://doi.org/10.1111/cobi.12872>

677 Wang, H.M. *et al.* (2019). Does the weighting of climate simulations result in a better quantification of  
678 hydrological impacts? *Hydrol. Earth Syst. Sci.* **23**, 4033–4050. [https://doi.org/10.5194/hess-23-](https://doi.org/10.5194/hess-23-4033-2019)  
679 [4033-2019](https://doi.org/10.5194/hess-23-4033-2019)

680 Willcock, S. *et al.* (2016). Do ecosystem service maps and models meet stakeholders’ needs? A  
681 preliminary survey across sub-Saharan Africa. *Ecosyst. Serv.* **18**, 110–117.  
682 <https://doi.org/10.1016/j.ecoser.2016.02.038>

683 Willcock, S. *et al.* (2019). A Continental-Scale Validation of Ecosystem Service  
684 Models. *Ecosystems* **22**, 1902–1917. <https://doi.org/10.1007/s10021-019-00380-y>

685 Willcock, S. *et al.* (2020). Ensembles of ecosystem service models can improve accuracy and indicate  
686 uncertainty. *Sci. Total Environ.* **747**, 141006. <https://doi.org/10.1016/j.scitotenv.2020.141006>

687 Wong, C.P. *et al.* (2014). Linking ecosystem characteristics to final ecosystem services for public  
688 policy. *Ecol. Lett.* **18**, 108–118. <https://doi.org/10.1111/ele.12389>

689 Zulian, G. *et al.* (2018). Practical application of spatial ecosystem service models to aid decision  
690 support. *Ecosyst. Serv.* **29**, 465–480. <https://doi.org/10.1016/j.ecoser.2017.11.005>  
691

**ANONYMISED MANUSCRIPT**

**Weighted Ensembles Reduce Uncertainty in Ecosystem Service Modelling**  
**Reducing Uncertainty in Ecosystem Service Modelling through Weighted Ensembles**

**Highlights:**

- Ensembles of models are used for other disciplines but not ecosystem services ~~(ES)~~
- ~~It is not known how~~ How best to combine ecosystem service ~~ES~~ models into an ensemble ~~is unknown~~
- We test ten contrasting ensemble approaches
- Ensembles had up to 27% higher accuracy than a randomly selected individual model
- Weighted ensembles provided better predictions

**Abstract:** (150 words)

Over the last decade many ecosystem service (ES) models have been developed to inform sustainable land and water use planning. However, uncertainty in the predictions of any single model in any specific situation can undermine their utility for decision-making. One solution is creating ensemble predictions, which potentially increase accuracy, but how best to create ES ensembles to reduce uncertainty is unknown and untested. Using ten models for carbon storage and nine for water supply, we tested a series of ensemble approaches against measured validation data in the UK. Ensembles had at minimum a 5-17% higher accuracy than a randomly selected individual model and, in general, ensembles weighted for among model consensus provided better predictions than unweighted ensembles. To support robust decision-making for sustainable development and reducing uncertainty around these decisions, our analysis suggests various ensemble methods ~~can~~ should be applied depending on data quality, for example if validation data are available.

**Graphical Abstract:**

Accuracy compared to mean ensemble	
<b>Individual Models</b>	Mean accuracy is always worse
<b>Unweighted Ensembles</b>	
Mean Ensemble	Reference ensemble type: up to 19% better than individual models
Median Ensemble	Mostly better, rarely worse
<b>Untrained Weighted Ensembles</b>	
Deterministic Consensus	Sometimes better, sometimes worse
Iterated Consensus	Mostly better, rarely worse
Attribute based	Sometimes better, sometimes worse
<b>Trained Weighted Ensembles</b>	
Accuracy weighted	Mostly better, rarely worse
Regressed consensus	Mostly better, never worse

29 **Keywords:** Carbon; Committee averaging; Prediction Error; Accuracy; United Kingdom; Validation;  
30 Water supply; Weighted averaging

31  
32 **Video Summary:** (see attached file)

### 33 1. Introduction

34 If the United Nations' sustainable development goals (SDG) are to be achieved worldwide (Griggs *et al.*  
35 2013), it is vital to understand and manage "nature's contributions to people" (termed ecosystem services;  
36 ES; Pascual *et al.* 2017). The empirical data needed to quantify ES are sparse in many parts of the world  
37 (Suich *et al.* 2015; Willcock *et al.* 2016), which is problematic as ES need to be accurately assessed and  
38 mapped to be incorporated in policy making and planning decisions (UKNEA 2011; de Groot *et al.* 2012).  
39 Such decisions require assessment of multiple ES, and the synergies and trade-offs among these ES, in order  
40 to estimate potential effects of land/water use change or other impacts (Willcock *et al.* 2016). Spatially-  
41 explicit models produce maps of estimated ES – typically based on globally available datasets of land cover  
42 combined with other predictor variables – and so can provide credible information of the spatial distributions  
43 of multiple ES, particularly where empirical data are lacking (Malinga *et al.* 2015; Costanza *et al.* 2017).  
44

45 Over the last 10 years, many ES models have been developed, by different teams, often using dissimilar  
46 approaches, and with little reference to the other models (Bagstad *et al.* 2013; Ochoa & Urbina-Cardona  
47 2017). For example, carbon stocks for climate change mitigation can be modelled by 'look-up tables'  
48 relating land cover to stocks, by deterministic statistical inference, or by simulating complex processes  
49 (Willcock *et al.* 2019). However, most applications of ES models rely on only a single model for each ES  
50 (Englund *et al.* 2017; Bryant *et al.* 2018). Furthermore, while models can only approximate reality, few  
51 applications explicitly validate ES models against independent datasets (Chaplin-Kramer *et al.* 2019),  
52 although there are notable exceptions (Redhead *et al.* 2016; Sharps *et al.* 2017; Willcock *et al.* 2019). This  
53 is a particular issue as the results of location-specific validation (*e.g.* that performed during model  
54 development) may not be transferable to new locations (Redhead *et al.* 2016), or up-scalable to the regional  
55 and national extents over which ES model outputs are required to achieve the SDG (Willcock *et al.* 2016;  
56 Willcock *et al.* 2019). From a user and stakeholder perspective, not knowing the accuracy of the available  
57 ES models for the region of interest typically leads to either selection of a single suboptimal model – at  
58 worst leading to perverse decision-making – or a reluctance to use ES models altogether, causing an  
59 implementation gap between research, incorporation into policy and subsequent decision-making (Wong *et al.*  
60 2014; Willcock *et al.* 2016).  
61

62 Despite claims for predictive superiority of certain modelling techniques and platforms, independent  
63 evaluations have been unable to demonstrate the pre-eminence of any single approach. In fact, while more  
64 complex models on average perform better in terms of fit to validation data, the best-fit model varies  
65 regionally and often according to the validation data used (Sharps *et al.* 2017; Willcock *et al.* 2019; Willcock  
66 *et al.* 2020). So, if no single ES model is always the most accurate, how should a suitable approach be  
67 selected?  
68

69 Across the sciences, one solution to address uncertainty surrounding the accuracy of any single model is to  
70 use an ensemble of models (Araújo & New 2007; Willcock *et al.* 2020) – using individual models as  
71 replicates with different input parameters and boundary conditions (Araújo & New 2007; Dormann *et al.*  
72 2018). Variation among models in their assumptions and formats can result in large differences in  
73 predictions, in terms of predicted values and how they vary over space, especially when there is uncertainty  
74 as to the state and processes of the system being modelled (van Soesbergen & Mulligan 2018; Willcock *et al.*  
75 2019). Ensembles of models are hypothesised to have enhanced accuracy over individual models due to  
76 fewer overall errors in prediction by reducing the influence of idiosyncratic outcomes from single models  
77 (Araújo & New 2007; Dormann *et al.* 2018). Individual models rarely capture all potentially relevant  
78 processes or are often tuned to particular ecosystem characteristics. A combination of models might provide  
79 a more comprehensive coverage of processes and their forms, and avoids the chance of (unknowingly)

80 selecting a model with a high prediction error at the location and scale of interest for a particular study  
81 (Willcock *et al.* 2020).

82  
83 Model ensembles are common in other disciplines – *e.g.* in niche modelling (Araújo & New 2007,  
84 Grenouillet *et al.* 2011), agroecology (Refsgaard *et al.* 2014), hydrology and water resources management  
85 (Wang *et al.* 2019; He *et al.* 2021), and climate and weather modelling (Knutti *et al.* 2013), as well as market  
86 forecasting (He *et al.* 2012). However, ensembles have been largely neglected in ES studies (Bryant *et al.*  
87 2018). The only current exception is the simplest ensemble approach (*i.e.* ‘committee averaging’ – taking  
88 the unweighted mean of a group of individual models per location –) which was applied to ES models in  
89 Sub-Saharan Africa, and gave higher accuracy in terms of fit to validation data (Willcock *et al.* 2020).  
90 Approaches that use more information might yield even more accurate estimates. Thus, here we explore the  
91 outstanding question of ‘what are the best ways to build ES model ensembles to realise the benefits such  
92 ensembles can bring to sustainability science?’

93  
94 Approaches to building model ensembles vary across disciplines, ranging from committee averaging  
95 (Marmion *et al.* 2009; Grenouillet *et al.* 2011) to complex Bayesian algorithms (Tebaldi & Knutti 2007).  
96 For example, species distribution models are generally deterministic statistical models; their fit to the data  
97 is often assessed with an accuracy metric and so ensembles are generally created using weighted averaging  
98 based on accuracy (Araújo & New 2007). By contrast, climate models are often treated as equal replicates  
99 with identical weights when making an ensemble (Tebaldi & Knutti 2007; Grenouillet *et al.* 2011) – we  
100 refer to such ensembles as ‘unweighted’. This difference may stem from the availability of suitable  
101 validation data, as well as different traditions. For example in species distribution models, biodiversity data  
102 are readily available and are used to train through cross-validation (Araújo & New 2007), whereas validation  
103 data on future climates obviously do not exist – although cross-validation against historic climate data is  
104 possible.

105  
106 As well as varying considerably in their underlying method, ES models often differ in the forms of their  
107 outputs (~~*e.g.* summed monetary value of the ES (de Groot *et al.* 2012) vs. specific biophysical predictions~~),  
108 even when modelling the same ES (~~*e.g.* summed monetary value of the ES (de Groot *et al.* 2012) vs. specific~~  
109 ~~biophysical predictions~~). By contrast, climate models generally have very similar forms of outputs. An  
110 important knowledge gap is therefore how to combine distinct ES model outputs as complementary inputs  
111 to provide a reliable ensemble. Outputs from different ES models can have different units and it is  
112 challenging to decide the relative weighting to place on each model. ~~M, with potentially different units, to~~  
113 provide reliable ensemble products using different model approaches as complementary inputs, and the  
114 potential role of weighting doing so. Since models for a particular ES often have different structures, may  
115 include different processes, or may represent the same processes in different ways (Ochoa & Urbina-  
116 Cardona 2017). As a result, the different ES models, they will most likely not have equal accuracy, and so  
117 prediction errors (*i.e.* bias) will may not be normally distributed among models (Dormann *et al.* 2018). If  
118 ES models had equal overall accuracies, unweighted averaging may provide a smoothing effect, reducing  
119 the impact of idiosyncratic outputs (*e.g.* at specific locations) of any particular model to reveal useful signals  
120 (Araújo & New 2007, Knutti *et al.* 2013; Diengdoh *et al.* 2020). In cases of varying overall accuracy,  
121 appropriate weighting of outputs based on model accuracy – *i.e.* models having unequal assigned weights –  
122 might re-adjust the distribution of prediction errors, and so improve the accuracy of the resulting ensemble  
123 (Refsgaard 2014; Dormann *et al.* 2018; Liu *et al.* 2020).

124  
125 However for ES, the lack of *a priori* validation data in many cases means that the distributions of accuracy  
126 among ES models are unknown. Furthermore, given that inferences about model accuracy at one location  
127 may not be transferable to others (Willcock *et al.* 2019), weighting using validation results from a separate  
128 study may not improve outcomes. Therefore where validation data are not available, the consensus among  
129 models could be used to weight their individual contribution to the ensemble value (Marmion *et al.* 2009;  
130 Grenouillet *et al.* 2011). This approach follows the logic that models whose output values are more different  
131 to those of the other models (*i.e.* are more distinct) are more likely to be incorrect. Therefore, weighting by  
132 consensus reduces the impact of outputs from more idiosyncratic models (*i.e.* those with extreme values,

133 outliers or badly comparable processes) by comparison with the other models (Araújo & New 2007;  
 134 Dormann *et al.* 2018), but does not exclude their information fully. The opposite may also be true – *i.e.*  
 135 more distinct models are more accurate – for example in cases where more similar models have common  
 136 inaccuracies.

137  
 138 Here, we implement 10 alternative ensemble methods, restricting ourselves to methods feasible for a wide  
 139 range of users, to evaluate whether weighting provides higher accuracy and if so which type of method  
 140 produces the most accurate predictions against validation data. We focus on two services, water supply and  
 141 carbon storage, in the United Kingdom. To support decision-making, we map the results for potential further  
 142 use, which ~~will be made~~ available via <https://doi.org/10.5285/a9ae773d-b742-4d42-ae42-2b594bae5d38>.  
 143 ~~through eide.ac.uk/~~ We use post-processing – specifically normalisation and per area correction –  
 144 developed in earlier work (Willcock *et al.* 2019; Willcock *et al.* 2020) to make outputs among models  
 145 comparable.

## 147 2. Methods

148 We developed and validated unweighted average and weighted average ensembles of models for a  
 149 provisioning service (water supply; subsequently referred to as ‘water’) and a regulating service  
 150 (aboveground carbon storage; subsequently referred to as ‘carbon’), for which there is both a variety of  
 151 models available (Bagstad *et al.* 2013; Ochoa & Urbina-Cardona 2017; Willcock *et al.* 2019) and the  
 152 presence of accessible validation data. We applied the models and ensemble methods in the United Kingdom  
 153 (UK), for which there is a large quantity of reliable validation data; allowing us to assess ensemble  
 154 accuracies. We compared accuracy (*i.e.* fit to validation data) of these individual models with those of the  
 155 ensembles generated from them via multiple approaches, assessed if weighted ensembles were an  
 156 improvement on the unweighted mean-averaged ensemble, and identified the methods of weighting  
 157 ensembles that gave the highest accuracy.

158  
 159 We modelled each ES at a 1 ha (100 × 100 m) resolution, and subsequently assessed performance of the  
 160 different ensemble approaches using weighting approaches we organised into three categories (Table 1):  
 161 deterministic consensus (*i.e.* always providing the same result), iterated consensus (*i.e.* using structured  
 162 trial-and-error approaches) and attribute-based (*e.g.* ~~spatial resolution~~ grain or distinctiveness). Finally, we  
 163 assessed the transferability of our UK results using independent data and models from a very different study  
 164 area – Sub-Saharan Africa (Willcock *et al.* 2019). We depict our overall process in Figure 1 in 7-steps. Our  
 165 calculations were performed using Matlab v7.14.0.739 and ArcMap 10.7.1, employing `ArcappPy` coding for  
 166 loops. Relevant codes can be found at [github.com/EnsemblesTypes](https://github.com/EnsemblesTypes), with flow among codes explained in  
 167 SI-1-3.

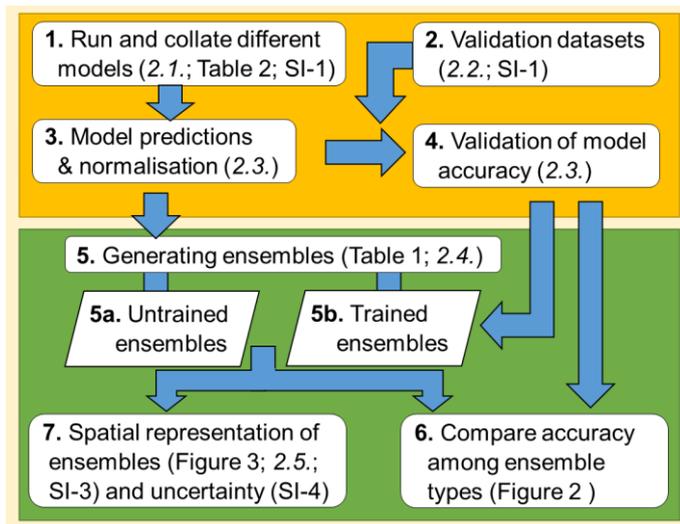
168  
 169 **Table 1. Approaches used to calculate accuracy (A) and ensembles (B).** Ensemble approaches were  
 170 applied to the outputs of ten models for carbon storage and nine for water supply (see Table 2). For weighted  
 171 averaging, the procedure is described, and where applicable the Matlab tools used are mentioned; similar  
 172 regression tools are available in most statistical packages (further explanation is provided in SI-1). Trained  
 173 weighting (En-9 & En-10) uses validation data, whereas untrained weighting (En-3 to En-8) does not. En-1  
 174 and En-2 are unweighted average ensemble approaches, and En-3 to En-10 are weighted average  
 175 approaches; the latter comprising *deterministic* (En-3 & En-4), *iterated* (En-5, En-6 & En-10) and *attribute*  
 176 *weighted* (En-7 to En-9) techniques. With  $\omega_i$ : weight for model  $i$ ;  $E_{(x)}$ : the value of the ensemble;  $V_{(x)}$ : the  
 177 normalised validation value;  $Y_{i(x)}$  and  $Y_{j(x)}$ : the normalised value of model  $i$  or comparator  $j$  respectively, all  
 178 for selected spatial point  $x$ ; ( $y \neq x$ ) denoting a split dataset;  $C_{(i,j)}$ : the correlation coefficient between model  
 179  $i$  and  $j$ ; with  $n$  the # models,  $m$  the # spatial data points;  $n^*$ : the # models in distinctiveness group  $g$  (see SI-  
 180 1 for distinctiveness grouping).

181

Approach	Description	Details & Matlab Tool
A. Accuracy approaches		

• Spearman $\rho$	Correlation coefficient between ranked variables $V$ and $T$ .	$T$ is either $Y_i$ or $E$ , depending on ensemble method	
• Inverse Deviance ( $D^{\downarrow}$ )	$D^{\downarrow} = 1 - \left( \frac{1}{m} \times \sum_x^m  X(x) - T(x)  \right)$	$T(x)$ is either $Y_{i(x)}$ or $\underline{E}_{i(x)}$	
<b>B. Ensemble approaches</b>			
<b>Unweighted Averaging:</b>			
En-1. Mean	$E_{(x)} = (\bar{Y}_i)_{(x)}$		
En-2. Median	$E_{(x)} = (\tilde{Y}_i)_{(x)}$	Hypothesised to perform better than mean for skewed distributions.	
<b>Untrained Weighted Ensembles: <math>E_{(x)} = \sum_i^n \left( \frac{\omega_i}{\sum_i^n \omega_i} \times Y_i \right)_{(x)}</math> with <math>\omega_i</math> following:</b>			
Deterministic consensus	En-3. PCA	$\omega_i =$ loadings of first Principal Component axis	Princomp-tool
	En-4. Correlation coefficients	$\omega_i = \frac{1}{n} \times \sum_j^n \frac{C_{(i,j)}}{\sqrt{C_{(i,i)} \times C_{(j,j)}}}$ , for all $j \in i$ with $C_{(i,j)} = \frac{1}{m-1} \times \sum_x^m \left( (Y_{i(x)} - \bar{Y}_i) \times (Y_{j(x)} - \bar{Y}_j) \right)$	
Iterated consensus	En-5. Regression to the median	$\tilde{Y}_{(x)} \sim (\sum_i^n \omega_i Y_i)_{(x)}$	nlmefit-tool, maximising Log Likelihood
	En-6. Exhaustive leave-one-out cross-validation <sup>3</sup>	$Y_{j(x)} \sim \sum_{i \neq j}^n \omega_i Y_{i(x)}$ , for all $j \in i$ subsequently: $\omega_i = \frac{1}{n} \times \sum_i^n \left( \left( \frac{1}{n-1} \right) \times \sum_{i \neq j}^n \omega_{ij} \right)$	nlmefit-tool, maximising Log Likelihood
Attribute-based	En-7. Upweighted <small>small finer spatial resolution</small> grains	$\omega_i = \frac{1}{\log_{10}(\text{spatial resolution} \neq \text{rain})}$	<small>Finer spatial resolution</small> Grain: the smaller grid size in 1-dimensional meters (e.g. 25 m)
	En-8. Attribute weighting: distinctiveness	$\omega_i = \left( \frac{n^g}{n} \right)$ when upweighted with $n^g = i \in g$ $\omega_i = \left( \frac{n}{n^g} \right)$ when downweighted with $n^g = i \in g$	
<b>Trained Weighted Ensembles: <math>\omega</math>-transfer via jack-knife training</b>			
Attribute-based	En-9. Accuracy-weighted	$\omega_i = A_i$ , with $A_i (V_{(y \neq x)}, Y_{(y \neq x)})$	With $A$ , either Spearman $\rho$ or $D^{\downarrow}$ accuracy
Iterated consensus	En-10. Log-likelihood regressions	$V_{(y \neq x)} \sim (\sum_i^n \omega_i Y_i)_{(y \neq x)}$	Using nlmefit-tool, maximising Log Likelihood

182



183

184 **Figure 1.** Schematic representation of our ensemble analysis with arrows  
185 showing information flows. Numbers represent the steps with the method chapters  
186 indicated in italics, with respective detailing SIs; result figures are indicated.  
187 Parallelograms highlight the 10 ensembles approaches (Table 1), using models  
188 described in Table 2.

189

190

### *2.1. Run and collate different models (step 1)*

191

192

193

194

195

196

197

198

199

200

201

202

203

We used outputs from 10 models for above ground carbon stocks based on per grid<sub>cell</sub> estimates, and outputs from nine models for annual water supply which provided accumulated flow estimates through specific pour points, either directly or through summation of run-off estimates per grid<sub>cell</sub>. We list these models in Table 2, including their output grid sizes (*spatial resolution<sub>grain</sub>*); we refer to SI-1-1 for full details, scales and supporting data. Acknowledging that model outputs have different units and sometimes model different constructs, we refer further to them in the general terms of carbon and water supply. Adhering to the aim of this paper, we do not compare individual model outputs, but focus on ensemble methods. All model outputs were set to the British National Grid transverse Mercator projection (EPSG 27700) with a 0.9996 scale factor and units in metres. Not all models covered the whole of the UK, *e.g.* some excluded Northern Ireland or Scotland (see SI-1-1). Where applicable we corrected for this by using a standard error of means as  $\left(\frac{\sigma(x)}{\sqrt{n(x)}}\right)$ , instead of standard deviation ( $\sigma$ ), with  $n$  the number of models per grid cell  $x$ . We collated models for this study according to their availability and to reflect different approaches to modelling ES.

204  
205  
206

**Table 2. Models and existing outputs used.** Full details, input data, post processing descriptions, and coverage are provided in SI-1-1. Model names are shown as acronyms and in full.

Model	Description	Grid size ( <u>spatial resolution</u> grain)	Model Type <sup>16</sup>
InVest v3.7.0 <sup>1†</sup> (Integrated Valuation of Ecosystem Services and Trade-offs)	Carbon module: above ground stocks	25 × 25 meters	Look-up table
	Water yield module: run-off per cell		Process
LPJ-GUESS <sup>2,3†</sup> (Lund-Potsdam-Jena General Ecosystem Simulator)	Vegetation biomass stocks per cell, mean for years 2009-2018	0.5° (≈ 46 × 46 km)	Process
	Water run-off per cell, mean for years 2009-2018		
LUCI <sup>4†</sup> (Land Utilisation Capability Indicator)	Above ground carbon stocks	10 × 10 meters	Look-up table
	Accumulated water run-off	5 × 5 meters	Process
\$-benefit transfer using The Economics of Ecosystems and Biodiversity database <sup>5,6†</sup>	Above ground carbon stock as monetary value	25 × 25 meters	Look-up table
	Water run-off as monetary value per cell		
Aqueduct v2.1 Total Blue Water <sup>7§</sup>	Accumulated water run-off	138 flow areas	Deterministic
ARIES k-Explorer <sup>8‡</sup> (Artificial Intelligence for Environment & Sustainability)	Joined above and below ground carbon stocks	1-hectare	Look-up table
Barredo <i>et al.</i> (2012) <sup>§</sup>	A European map of above ground biomass stocks	1 km <sup>2</sup>	Look-up table
Copernicus, Tree Cover Density <sup>9§</sup>	Proxy for carbon: tree Cover Density 2015 from MODIS satellite imagery.	20 × 20 meters	Deterministic
DECIPHeR <sup>10§</sup> (Dynamic fluxEs and ConnectIvity for Predictions of HydRology)	Accumulated water run-off through NRFA delineated catchment outlets, mean for years 1995-2015	387 catchments in common with validation	Process
Grid-to-Grid <sup>11§</sup>	Accumulated water run-off, mean for years 1995-2015	1 km <sup>2</sup>	Process
Henrys <i>et al.</i> (2016) <sup>§</sup>	Above ground carbon stocks	1 km <sup>2</sup>	Look-up table
Kindermann <i>et al.</i> (2008) <sup>§</sup>	A global map of above ground forest biomass stocks	1 hectare	Deterministic
National Forest Inventory (2018) <sup>12†</sup>	Woodland Land Cover Map <sup>15</sup> with above ground carbon stocks based on added Look-up table (Table. SI-1-4)	20 × 20 meters	Look-up table
Scholes Growth Days <sup>13,14†</sup>	Proxy for water run off per cell: # Days precipitation exceeds evapotranspiration	1 km <sup>2</sup>	Deterministic
WaterWorld v2 <sup>15‡</sup>	Accumulated water run-off	0.0083° (≈ 1 km <sup>2</sup> )	Process

207  
208

<sup>†</sup>Output generated for this work; <sup>‡</sup>online tool; <sup>§</sup>existing dataset; <sup>1</sup>Kareiva *et al.* (2011); <sup>2</sup>Smith *et al.* (2014); <sup>3</sup>Ahlström *et al.* (2015); <sup>4</sup>Thomas *et al.* (2020); <sup>5</sup>de Groot *et al.* (2012); <sup>6</sup>Costanza *et al.* (2014); <sup>7</sup>Gassert *et al.* (2015) <sup>8</sup>Martínez-López *et al.* (2019); <sup>9</sup>[land.copernicus.eu/tree-cover-density/status-maps/2015](http://land.copernicus.eu/tree-cover-density/status-maps/2015); <sup>10</sup>Coxon *et al.* (2019a; 2019b);

Field Code Changed

209 <sup>11</sup>Bell *et al.* (2018a; 2018b); <sup>12</sup>Forestry Commission (2018); <sup>13</sup>Scholes (1998); <sup>14</sup>Willcock *et al.* (2019); <sup>15</sup>Mulligan (2013); <sup>16</sup>following Ding & Bullock (2018), Willcock *et al.*  
210 (2019).  
211

## 2.2. Validation datasets (step 2)

Our carbon stock validation dataset was provided by Forest Research and comprises species inventories in all forest estates in England and Scotland in 2019 ([data-forestry.opendata.arcgis.com/](https://data-forestry.opendata.arcgis.com/); density shown in Figure 3; locations in Figure SI-1-2). In 201,143 forest compartments of varying size (mean: 4.4 hectares, median 1.6 hectares,  $\pm 22.1$ ), tree species, stand age and thinning regime were recorded for three vegetation layers. For each compartment and layer therein, the unique combination of stand age, thinning regime and tree species of the inventory data was searched in the UK Carbon Code tables ([woodlandcarboncod.org.uk](https://woodlandcarboncod.org.uk)) and life-time accumulated biomass was converted to total standing carbon per hectare estimates per compartment, with the layers summed per compartment (SI-1-2). Subsequently, compartments were spatially joined into 2078 polygons of 'forest' that were separated if more than 25 meters distance from each other.

Our water supply validation dataset comprised 519 hydrometric gauging stations from the National River Flow Archive of the UK (NRFA; [nrfa.ceh.ac.uk](https://nrfa.ceh.ac.uk)), with associated catchments representing a variety of sizes distributed across the whole of the UK (Figure 3). From the 1598 potential catchments in NRFA, we selected those that were  $>100 \text{ km}^2$  to get a robust mean run-off from the catchments. In cases where multiple gauging stations were found along the same river, based on name, only the largest was chosen to avoid pseudoreplication. An additional set of 41 Welsh catchments was included which did not meet this size criterion. Wales contains mainly small catchments due its geography – mountain ranges close to the sea – and so we selected catchments  $>25 \text{ km}^2$  to avoid this part of the UK being underrepresented. The data were polygons encompassing these catchments. Details are provided in SI-1-2.

## 2.3. Model predictions, normalisation (step 3) and validation of model accuracy (step 4)

For each individual model, predictions were obtained for each polygon in the validation dataset using the ArcGIS spatial analyst Zonal tool with a forced 2.5 m grid size environmental setting to minimise edge effects; *i.e.* all predicted values were obtained by resampling into  $2.5 \times 2.5 \text{ m}$  grid cells. In most cases the modelled value per polygon was obtained by taking the sum of all constituent grid cell values, corrected for both actual grid size and the resampling to 2.5 m. In the case of accumulated flow models, we corrected for potential small scale differences in flow routing among these models by taking the maximum flow value within both a 2 km range of the NRFA reported location of the gauging station and the polygon associated with that gauging station.

To ensure comparability among model outputs, we standardised by normalising among the outputs for each individual model and for the validation data-sets. Prior to this step all outputs were area corrected as either mean carbon stock – or proxy thereof – per hectare or water supply per hectare of catchment (with accumulated run-off estimates post-processed to give net run-off per cell; SI-1-1). This normalisation followed Willcock *et al.* (2019), and allowed us to address differences in units among models (such as monetary benefit transfer vs. satellite-based tree cover densities or run-off, and equalised carbon and biomass). To avoid impacts of extreme values without eliminating such data-points, we employed a double-sided Winsorising protocol for normalisation (Willcock *et al.* 2019; Verhagen *et al.* 2017), using the values associated to the 2.5% and 97.5% percentiles of number of datapoints to define the 0 and 1 values (values below or above these percentiles became 0 or 1 respectively). This winsorising normalisation protocol assumes outlier data are valid, but skewed values, in our case mainly by per area averaging, and corrects for this by compressing the variance tails rather than trimming them (Keselman *et al.* 2008; Erceg & Miroseovich 2008). Hence, we trade-off an even data distribution over the full 0-1 normalised range against the chance of having a true far outlier maximum (see SI-5 for a full investigation into the impact of the Winsorising protocol over standard normalisation for the validation data distribution). For each model, normalisation was done prior to creating ensembles.

For validation, we employed two accuracy measures (Willock *et al.* 2019; Willock *et al.* 2020), which are related to different aims in modelling ES (Table 1):

- 263 1) Comparing the rank order of predicted and validation data using Spearman  $\rho$ . This is relevant where  
 264 modelling is used to discover, for example, the most important locations for delivering an ES, or  
 265 conversely, those areas whose development may have least impact on ES delivery.  
 266 2) Ascertaining the absolute difference of each modelled value from its validation value using the inverse  
 267 of the deviance ( $D^{\downarrow}$ ). This is relevant where modelled values are important, *e.g.* when testing where ES  
 268 levels exceed a minimum threshold. We used the inverse of the deviance so that, like  $\rho$ , a higher value  
 269 indicated greater accuracy.

270  
 271 *2.4. Generate ensembles (step 5) and compare accuracy among ensemble types (step 6)*

272 We tested whether model ensembles were more accurate than the individual constituent models and which  
 273 approaches for creating ensembles were the most accurate in terms of fit to validation data. We created  
 274 ensembles using a range of methods, from the simplest calculation of an average value of the models at each  
 275 location ('unweighted averaged ensembles', *e.g.* Marmion *et al.* 2009, Grenouillet *et al.* 2011) to ensembles  
 276 with the contributions from different models weighted unequally ('weighted ensembles'), following  
 277 Dormann *et al.* (2018) (Table 1; further explanation and a model flow are provided in SI-1-3). We used  
 278 relatively straightforward approaches that would be feasible for a wide community of scientists and  
 279 decision-makers, and avoided more complex mathematical and/or statistical techniques such as Bayesian  
 280 networks (Bryant *et al.* 2018), which would require detailed specialist knowledge. Weights over all models  
 281 were normalised to sum to 1. Together with normalisation of the ensemble outputs (see above), this assured  
 282 equal scaling among all models and ensembles.

283  
 284 For unweighted average ensembles, we calculated both the mean and the median of modelled values at each  
 285 location as alternative measures of the central tendency which are differently affected by skew in the data  
 286 (Table 1, En-1 & En-2).

287  
 288 For weighted ensembles we calculated:

289 
$$E_{(x)} = \sum_i^n \left( \frac{\omega_i}{\sum_i^n \omega_i} \times Y_i \right)_{(x)}$$
 with positive weights  $\omega_i$  for model  $i$  of validation polygon  $x$ , weights  $\omega_i$  are  
 290 normalised to sum to 1,  $Y$  the modelled values for  $i$  per polygon (step 3), and  $n$  the total number  
 291 of models per service.

292  
 293 To determine  $\omega_i$ , the weighting value for each model  $i$ , we employed a range of methods that can be broadly  
 294 categorised as two main types of ensemble approach (untrained and trained), with further subdivision as:  
 295 deterministic consensus, iterated consensus, and attribute-based. The ensembles are listed as equations in  
 296 Table 1 (see SI-1-3 for further details).

297 1) Untrained ensembles (En-3 to En-8) represent a situation in which there is no validation data. To generate  
 298 uncertainty estimates allowing statistical comparison with the models and among ensembles we jack-  
 299 knifed (Araújo & New 2007; Refsgaard *et al.* 2014) with 50% of the spatial data polygons for 250 runs,  
 300 *i.e.* every run contained a new selection of half the dataset. We tested three approaches to produce the  
 301 ensembles:

- 302 - *Deterministic consensus* among models can be calculated using several approaches, including the fit  
 303 to a common consensus axis such as from a Principal Components Analysis (Marmion *et al.* 2009;  
 304 Grenouillet *et al.* 2011) or weighting by correlation coefficients (En-3 & En-4; ensemble numbering  
 305 follows Table 1).
- 306 - *Iterative approaches* might more accurately quantify consensus among models through using  
 307 structured trial-and-error (Dormann *et al.* 2018; Tebaldi & Knutti 2007). We use two regression  
 308 techniques: between the individual models and the median (En-5) and leave-one-out cross-validation  
 309 (En-6) following the suggestion in Dormann *et al.* (2018).
- 310 - One might *a priori* place value on a particular model attribute and use this to create weights (Englund  
 311 *et al.* 2017; Willcock *et al.* 2019; Brun *et al.* 2020; En-7, En-8 & En-9). For example, one could up-  
 312 or down-weight more distinct model types through a binary matrix of differences (En-8 & En-9; SI-

313 1-4) in land cover map used, grid-size, measured or modelled climate, model extent, presence of  
314 time-series, time step-size and model type (*i.e.* look-up table, deterministic or process based).  
315 Alternatively models that run at coarser spatial resolutions are penalised (En-7): smaller grid sizes  
316 are deemed more useful for decision-making (Willcock *et al.* 2016).

317 2) Trained ensembles (En-9 & En-10), as often used for species distribution models (*e.g.* Refsgaard *et al.*  
318 2014; Elith *et al.* 2011), represent a situation in which validation data are available from a similar region  
319 or part of the study area and so cannot be used to directly validate or substitute for the models in the  
320 study area, but can be used to weight these models. Here,  $\omega_i$  was trained with the validation data on a  
321 jack-knifed 50% of the dataset to achieve maximum accuracy (En-10) and subsequently  $\omega_i$  was  
322 transferred to the other half of the dataset. We used 250 such jack-knife runs (see above), with the same  
323 selections as above. Moreover, we included weighting by individual model accuracy (Marmion *et al.*  
324 2009; Liu *et al.* 2020) using the same jack-knife approach (En-9).

325  
326 After creating the ensembles, their accuracy was assessed following step 4 using the two measures (see 2.3):  
327 Spearman  $\rho$  and the inverse of the deviance ( $D^{\dagger}$ ). We assessed any improvement over the unweighted mean-  
328 averaged ensemble as the reference with pairwise t-tests against the null hypothesis of equal accuracy  
329 (Matlab *ttest*-tool). A similar analysis against the median-averaged ensemble as reference can be found in  
330 SI-2. To avoid spurious findings of significance through having a large number of replicates, we assessed  
331 improvement using bootstrapped tranches of 50 runs each with 250 replicates, and averaging the P-values.  
332 Since we used the same statistical test 12-times per service per accuracy estimate, we employed a full  
333 conservative Bonferroni correction; ( $\alpha = 0.05/12$ ) on the resulting average P-values. To compare the  
334 ensembles with the individual models we calculated per replicate the mean difference in accuracy among

335 all models ( $A_i$ ) against accuracy of an ensemble ( $A_E$ ) following:  $\left( \left( \sum_i^n \left( \frac{A_E}{A_i} - 1 \right) \right) \times \frac{1}{n} \right)$ , with n the number  
336 models and  $i$  an individual model.

337  
338 Steps 5 and 6 were repeated using independent data and models from a different study area (sub-Saharan  
339 Africa; Willcock *et al.* 2019) to investigate the transferability of the results presented here (Figure SI-2-2).  
340

### 341 2.5. Spatial representation of ensembles and uncertainty (step 7)

342 To better support decision-making, we mapped our ES ensembles for the UK. For all the water ensembles,  
343 the mean normalised value across jack-knifed ensemble predictions per ensemble method were mapped as  
344 catchment polygons (step 5, N = 519). For all carbon ensembles we mapped as 1 km<sup>2</sup> grid cells. Here, for  
345 each ensemble approach, the estimated weights as calculated for the validation polygons – mean averaged  
346 among jack-knife runs – were transferred to the full area, with the result aggregated to a 1 km<sup>2</sup> resolution  
347 based on the mean value among 1 hectare grid\_cells. In total, this carbon dataset has 253,802 cells that  
348 (partially) contain non-sea land cover. We transferred the weights calculated for the forests since running  
349 cross-validation approaches on over 250K data points would extremely time consuming to compute.  
350 However, since our validation data are only from forests/woodlands, we are aware of introducing a potential  
351 bias that could skew non-forested areas to lower values. Furthermore, we generated UK-scale maps of  
352 spatial variation in the differences among the untrained ensemble approaches, by calculating the standard  
353 error of the mean (SEM) among these spatial outputs. These maps ~~will be made~~ are freely available online  
354 (<https://doi.org/10.5285/a9ac773d-b742-4d42-ae42-2b594bae5d38> ~~through [ecide.ac.uk](mailto:ecide.ac.uk)~~), and spatial  
355 patterns of uncertainty are discussed in SI-4.

356  
357

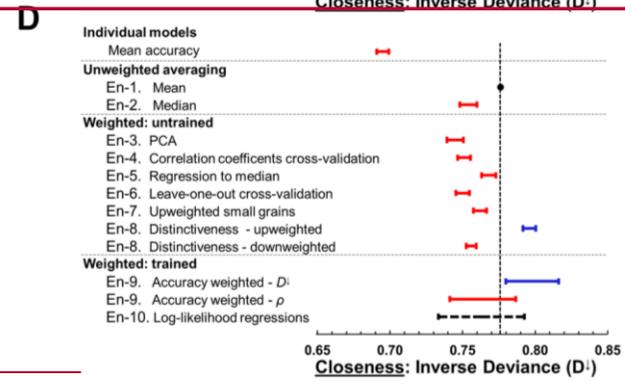
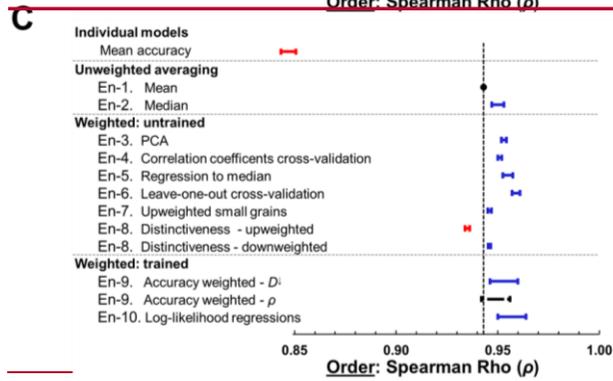
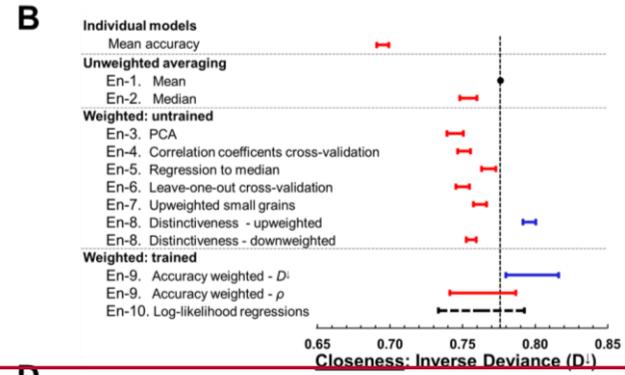
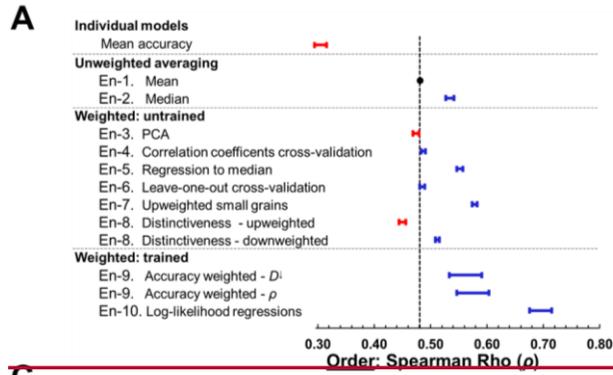
## 358 3. Results

359  
360

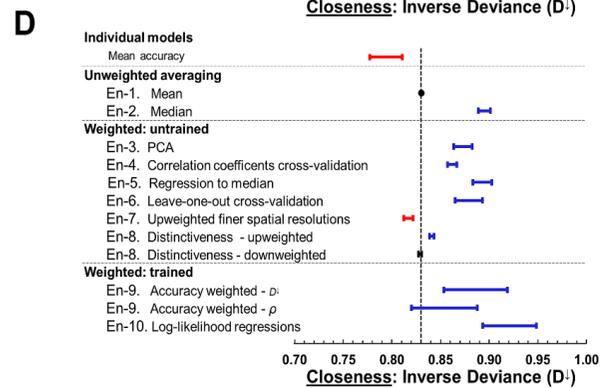
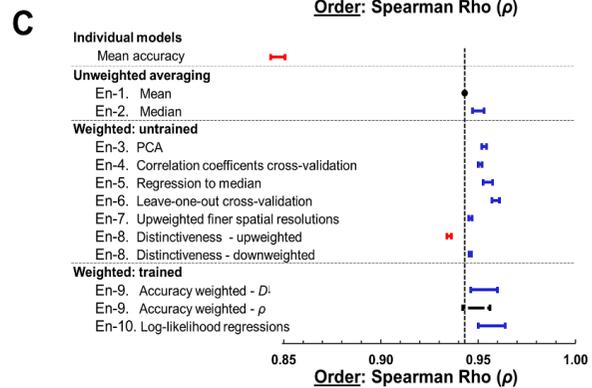
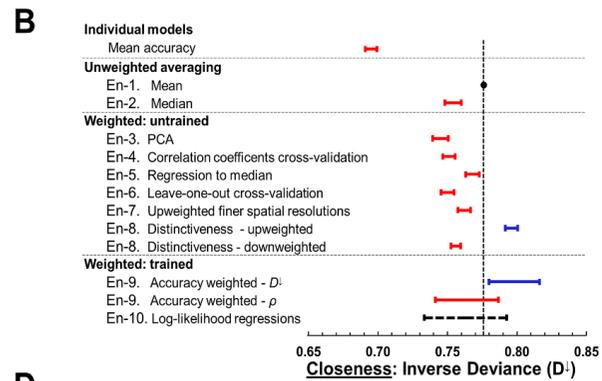
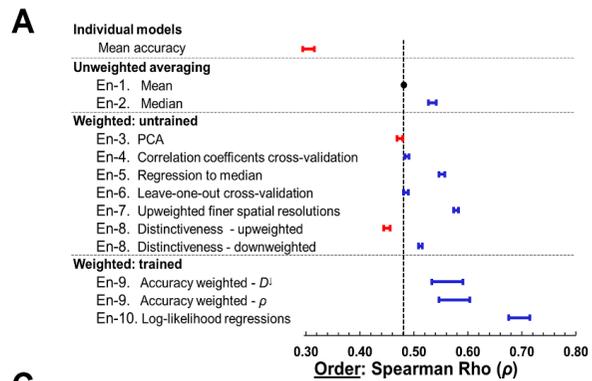
### 361 3.1. Ensembles are more accurate than individual models

362 The average accuracy of individual models, represented by the mean of accuracy values taken across all  
363 models, was lower than that for any of the ensembles we created. The accuracy of the unweighted averaged  
364 ensembles (of modelled values at each location, *e.g.* ‘mean ensemble’) was appreciably higher than the

363 mean value for accuracy of the individual models for both carbon and water: 19%  $\pm$ 1.1% [sd] for  $\rho$  and  
364 12.1%  $\pm$ 0.5% for  $D^1$  improvement in fit to the validation data for carbon and 5.7%  $\pm$ 0.4% for  $\rho$  and 9.5%  
365  $\pm$ 1.7% for  $D^1$  for water (Figure 2). Untrained weighted ensembles showed large improvements – for most,  
366 larger than the unweighted ensembles – over the mean accuracy of the individual models of 17% to 27%  
367 ( $\rho$ ) and 7.6% to 15% ( $D^1$ ) for carbon (Figure 2A and B), and 5.3% to 6.5% ( $\rho$ ) and 7.7% to 18% ( $D^1$ ) for  
368 water (Figure 2C and D). In all cases, pairwise t-tests indicated highly significant differences between each  
369 ensemble and the mean value of accuracy of individual models (all  $P < 1E^{-10}$ ). Thus, creating an ensemble  
370 improves prediction accuracy against a randomly chosen individual model irrespective of the ensemble  
371 approach chosen.



Formatted: Centered



374 **Figure 2. Accuracy of above ground carbon stock ensembles (10 models; A and B), and of water supply ensembles (9 models; C and D) against validation**  
375 **data.** The mean of accuracy values across the containing models – *i.e.* a randomly chosen model– is provided for comparison. For detail on the different ensemble  
376 types see Table 1 and SI-1-3. We show the average accuracy of 250 bootstrap runs with 50% of the dataset. The vertical dashed line indicates the reference  
377 unweighted mean-averaged ensemble (black dot, ‘mean ensemble’). Error bars indicate the standard deviation among runs in terms of proportional difference  
378 to the mean ensemble, calculated per bootstrap run as the difference in accuracy to the mean ensemble divided by the accuracy of the mean ensemble. The  
379 coefficient of variation among bootstraps for the mean carbon ensemble was 4% and 1%, for  $\rho$  and  $D^{\downarrow}$  respectively, and 1 % and 2% for water (not shown). **Blue**  
380 coloured ensemble accuracies are significantly higher than the unweighted mean ensemble (Bonferroni corrected  $\alpha = (0.05/12)$ ); **Red** coloured bars are  
381 significantly lower; **Black** dashed bars are not significantly different to the mean ensemble.

382 3.2. *Weighted ensembles are more accurate than unweighted ensembles*

383 All weighted ensembles, whether trained or untrained, significantly outperformed the reference unweighted  
384 mean ensemble (Figure 2), with the exception of  $D^{\downarrow}$  for carbon. In all cases, pairwise t-tests indicated these  
385 differences were highly significant ( $P < 1E^{-10}$ ; see Figure SI-2-1 for similar analyses against the median-  
386 averaged ensemble).

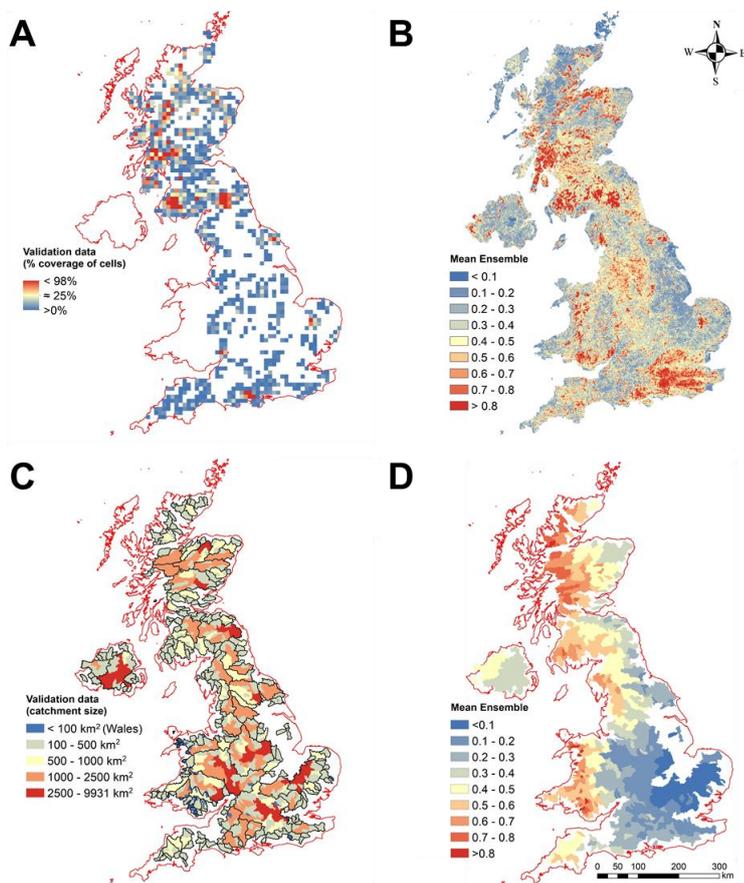
387  
388 For untrained weighted ensembles, prediction accuracy was elevated by up to  $4.8\% \pm 0.6\%$  for carbon  $\rho$   
389 (best: regression to median; Figure 2), with no improvement for carbon  $D^{\downarrow}$ , and  $0.8\% \pm 0.3\%$  and  $7.5\%$   
390  $\pm 1.1\%$  for water supply  $\rho$  and  $D^{\downarrow}$  respectively (regression to median; Figure 2). Conclusions as to the best  
391 model attributes to use for untrained weighting were dependent on the accuracy metric used ( $\rho$  or  $D^{\downarrow}$ ). By  
392 comparison to the unweighted mean ensembles, upweighting ~~smaller-grained~~ model outputs with finer  
393 spatial resolution improved  $\rho$  by up to  $6.6\% \pm 0.5\%$  and  $0.2\% \pm 0.1\%$  for carbon and water respectively but  
394 contrastingly decreased  $D^{\downarrow}$ . Upweighting more distinctive models was positive for  $D^{\downarrow}$  with  $2.5\% \pm 0.4\%$  and  
395  $1.3\% \pm 0.3\%$  greater accuracy compared to the unweighted mean ensemble for carbon and water supply  
396 respectively, but was negative for  $\rho$ . In summary, creating untrained weighted ensembles through iterative  
397 approaches was overall the most robust – particularly regression to the median (Table 1: En-5), showing  
398 greater accuracy than the unweighted mean-averaged ensembles in 3 out of 4 of our tests, and lower accuracy  
399 in 1 (Figure 2).

400  
401 For trained weighting ensembles, using an iterative log-likelihood regression approach (Table 1: En-10) to  
402 establish weights elevated prediction accuracy compared to the unweighted mean ensemble by up to  $14.5\%$   
403  $\pm 2.6\%$  for carbon  $\rho$  (no improvement for carbon  $D^{\downarrow}$ ) and  $0.8\% \pm 0.7\%$  and  $11.1\% \pm 3.4\%$  for water supply  $\rho$   
404 and  $D^{\downarrow}$  respectively (Figure 2). Compared to such regressions, upweighting models with higher accuracy in  
405 the training set (accuracy-weighted ensembles; En-9; Figure 2) gave less improvement over the unweighted  
406 mean ensemble. Iteratively creating trained weighted ensembles using a log-likelihood regression approach  
407 (Table 1: En-10) was most robust – showing greater accuracy than the unweighted mean-averaged  
408 ensembles in 3 out of 4 of our tests, and is no worse in 1 (Figure 2).

409  
410 The reference unweighted mean ensembles for carbon and water are mapped for the UK in Figure 3. Maps  
411 for all other ensembles can be found in SI-3 and uncertainty among models and ensembles in SI-4. In  
412 accordance with *a priori* predictions, the uncertainty associated with selecting a single model was several  
413 times greater than that associated with selecting any single ensemble method for both ES. For carbon, the  
414 standard error of the means (SEM) among individual models per 1 km<sup>2</sup> grid cell (SEM =  $9.0\% \pm 2.8\%$ , SI-  
415 4) was ca. 3.5-times larger than among ensembles (SEM =  $2.5\% \pm 1.1\%$ ). Similarly, the SEM among  
416 individual water models per watershed (SEM =  $7.8\% \pm 3.4\%$ , SI-4) was substantially greater than among  
417 ensembles (SEM =  $1.3\% \pm 0.7\%$ ). In SI-4 we investigate spatial drivers for this uncertainty, discussing these  
418 patterns at length.

419  
420 We validated the robustness of our results using independent data and models from a different area (Sub-  
421 Saharan Africa; Willcock *et al.* 2019), which gave similar results of weighted ensembles outperforming the  
422 reference mean ensemble (Figure SI-2-2).

423



424  
 425 **Figure 3. Spatial distribution of validation points and the reference mean ecosystem service value. A**  
 426 the Distribution of 2078 carbon validation forests as coverage of  $10 \times 10$  km cells – many individual forest  
 427 fragments would be too small to be clear at this scale, see SI SI-1-2 –, white cells are empty. **B** the reference  
 428 unweighted mean ensemble of carbon across 10 models, normalised on scale 0-1. **C** the 519 catchments  
 429 used for water validation and ensemble calculations coloured by their size – smaller watersheds that overlap  
 430 larger ones are displayed on top; lines show underlying largest catchment level. **D** the reference unweighted  
 431 mean ensemble of water supply across 9 models, normalised on scale 0-1. All maps here, in SI-3 (all  
 432 ensembles) and SI-4 (uncertainty) could support landscape decisions in the UK and ~~will be made~~are  
 433 available ~~through~~ via <https://doi.org/10.5285/a9ae773d-b742-4d42-ae42-2b594bae5d38eide.ac.uk/>.

434  
 435 **4. Discussion**

436 We have shown that predictions from ensembles of models have substantially higher accuracy than a  
 437 randomly selected single ES model, and especially that weighting approaches increase ensemble accuracy.  
 438 Finding increased performance through use of ensemble approaches is common in other fields. For example,  
 439 the increased accuracy of ensemble species distribution models ranges from 1-2% (Crossman *et al.* 2012;  
 440 Abrahms *et al.* 2019) to 12% (Grenouillet *et al.* 2011), although an increase is not universal (Hao *et al.*  
 441 2020). Similarly, 2% accuracy increases were found for market forecasting ensembles (He *et al.* 2012), and  
 442 neural network ensemble averaging resulted in up to 7% improvements in accuracy (Inoue & Narisha 2000).

443  
444 Specific to ES, unweighted averaged ensembles have been shown to be 5.0–6.1% more accurate than  
445 individual models (Willcock *et al.* 2020). Our improvements with ES ensembles are at minimum 5%-17%,  
446 suggesting substantial differences among models in their adequacy (Dormann *et al.* 2018), but also that  
447 ensemble approaches that use more information offer greater increases in accuracy. We found that taking  
448 the median generally outperforms a mean ensemble, probably because the latter is more influenced by  
449 outliers. Our results provide evidence that weighted ES ensembles created using consensus techniques  
450 produce more accurate outputs than unweighted ensembles. This finding is supported by our additional  
451 analysis using independent models and data from Sub-Saharan Africa (in a biome with very different  
452 climatic and soil characteristics; SI-2), suggesting our findings may be generalisable, although investigating  
453 this specifically (e.g., for different ES, regions and validation datasets) is an important avenue for future  
454 research.

455  
456 Predictions from models, including those from ES models, are all potentially biased in direction and amount  
457 because of their underlying assumptions. These biases could differ among models due to their specific  
458 construction. Therefore, models are likely to differ in their accuracy when compared to reality (Dormann *et al.*  
459 *et al.* 2018). The improvement in accuracy when using ensembles, as we have shown here, is referred to as a  
460 ‘portfolio effect’ by which a (weighted) combination of replications of possible states of a system suppresses  
461 idiosyncratic differences and provides a more reliable average estimate (Thibaut & Connolly 2013;  
462 Dormann *et al.* 2018; Lewis *et al.* 2021). However, this effect is lessened if models share similar  
463 assumptions and, therefore, concomitant biases – highlighting the importance of including multiple model  
464 outputs (Ding & Bullock 2018) and, where data are available, model validation (Willcock *et al.* 2019). In  
465 particular, the use of models not usually packaged as ES models – such as LPJ-GUESS – might help with  
466 increasing the variety of inputs for ensembles. If some models systematically overestimate and other models  
467 underestimate, averaging delivers smaller prediction errors when models are weighted (Dormann *et al.*  
468 2018). Hence, the resulting weighted ensemble is more accurate than most individual models and  
469 unweighted approaches (Marmion *et al.* 2009, Grenouillet *et al.* 2011); see Dormann *et al.* (2018) for  
470 theoretical explorations.

471  
472 We have shown the general potential of weighting to re-balance the contribution of different ES models,  
473 but also find that some weighting approaches seem more suitable. Specifically, structured trial-and-error  
474 iterative approaches may more accurately maximise consensus among models than deterministic approaches  
475 (Dormann *et al.* 2018; Gobeyn *et al.* 2019). The PCA and correlation coefficient approaches (Table 1: En-  
476 3 & En-4) deterministically assess consensus among individual models. By contrast, regression to the  
477 median, leave-one-out cross validation, and log-likelihood approaches (Table 1: En-5, En-6, En-10) are  
478 examples of iterative processes that optimise for the highest level of consensus in full parameter space  
479 (Dormann *et al.* 2018). Attribute-based approaches as used by Masson & Knutti (2011) and Willcock *et al.*  
480 (2019) (e.g. weighting by model distinctiveness or grid size; Table 1: En-7 and En-8) produce conflicting  
481 results. Model attributes such as these may not correctly describe why model outputs vary, or capture their  
482 complexity (Willcock *et al.* 2019; Brun *et al.* 2020) and so weighting by among-model agreement produces  
483 more accurate ensemble outputs. One might expect accuracy-weighted ensembles (Table 1: En-9) to  
484 perform best. However, model accuracy can be location specific and poorly transferable elsewhere – even  
485 with similar model accuracy, some grid cells may be well represented by some models and less by others  
486 (Graham *et al.* 2008; Marmion *et al.* 2009; Zulian *et al.* 2018). As a result accuracy-derived weights show  
487 high uncertainty in areas where training data were not available (i.e. non-forested areas; SI-4), likely because  
488 of over-fitting to areas with available data (i.e. forests/woodlands) producing correlative patterns that  
489 explain other areas less wellAs a result accuracy-derived weights show high uncertainty in areas where  
490 training data were not available (SI 4), likely because of over fitting to woodland areas. In SI-4, we  
491 investigated environmental and spatial drivers of uncertainty among predictions. Broadly, these  
492 supplementary results show that carbon models and ES ensembles are less accurate in urban areas. We also  
493 find that ensembles for water are less accurate in areas of high rainfall, seasonality and rugosity (see SI-4

Formatted: Font: (Default) + Headings CS (Times New Roman)

494 for full details). That said, as uncertainty among ES ensembles is almost 4-times lower than among  
495 individual models, this suggests less need to make the ‘right choice’ of method when selecting an ensemble  
496 approach. Thus, although there is some chance of picking a superior individual model (Willcock *et al.* 2018),  
497 the risk of a sub-optimal prediction is substantially lowered by applying any ensemble method and this risk  
498 is further reduced when a weighted ensemble is used.

499  
500 Our results should serve as a ‘call to arms’ for ES researchers and practitioners to increasingly use ensembles  
501 of models to support decision-making for sustainability. Using an individual ES model is fraught with  
502 concerns as *a priori* it is not known which is the most accurate and choosing only one model can, at worst,  
503 result in perverse decisions (Willcock *et al.* 2019). Deriving decisions from an ensemble of ES models  
504 provides an improvement over using one model for any location (which may be large or small, depending  
505 on the local context and the models used), but also more consistency over space, as model accuracy varies  
506 spatially (see results in SI-4). Therefore, using ensemble approaches, and especially weighted ensembles,  
507 would increase credibility and so help reduce the implementation gap between research and policy- and  
508 decision-making (Wong *et al.* 2014; Willcock *et al.* 2016). We acknowledge the lack of standardised metrics  
509 across models and limited computational and financial resources that could restrict the uptake of ensembles  
510 – indeed, many practitioners only run a single model. However, given the errors associated with single  
511 models (this paper; Willcock *et al.* 2020; Eigenbrod *et al.* 2010), we argue that a single model is inadequate,  
512 although more complex models are sometimes more accurate (Willcock *et al.* 2019). The most complex (a  
513 priori best) ES models require substantial inputs (i.e. data, computational power, subscription fees, and staff  
514 time), and so running multiple models – whilst requiring additional resources – results in a large gain per  
515 extra unit resource. For example, as even untrained weighted ensembles developed using iterative  
516 approaches (e.g. regression to the median, leave-one-out cross validation) enable a 3-fold reduction in  
517 variation, such an ensemble approach seems a reasonable minimum standard for ES modelling – striking  
518 the right balance between feasibility and robustness (Willcock *et al.* 2016). Whilst such ensembles will be  
519 outperformed by the best-performing individual models, these cannot be identified without running multiple  
520 models – a ‘Catch-22’ (Willcock *et al.* 2019). Thus, we recommend that multiple models be developed for  
521 ES where they are lacking (e.g. cultural services; Martínez-Harms and Balvanera, 2012; Wong *et al.* 2014),  
522 and that those with access to sufficient resources to run multiple models ensure the ensemble outputs are  
523 freely available, making the use of these ensembles more feasible and accessible for all (Willcock *et al.*  
524 2020).

## 526 5. Conclusion

527 We show that in situations with no *a priori* validation evidence guiding model selection, predictions from  
528 ensembles of models have a higher accuracy than selecting an individual model by chance. Weighted  
529 averaging further improves accuracy, suppressing idiosyncratic differences through producing consensus  
530 (Araújo & New 2007; Dormann *et al.* 2018). Doing so not only elevates accuracy but substantially decreases  
531 uncertainty among ensemble approaches compared to uncertainty among models, a further indication of  
532 increased fit to reality (Chaplin-Kramer *et al.* 2019; Willcock *et al.* 2020). In summary, even if a less  
533 accurate ensemble weighting approach is used, one would on average have lower uncertainty than selecting  
534 an individual model by chance. Thus, particularly when validation data are not available, we recommend  
535 the use of weighted ensembles in ES research to substantially reduce uncertainty and to support robust  
536 decision-making for sustainable development.

537

## 538 References

- 539 Abrahms, B. *et al.* (2019). Dynamic ensemble models to predict distributions and anthropogenic risk  
540 exposure for highly mobile species. *Divers. Distrib.* **25**, 11821193.  
541 <https://doi.org/10.1111/ddi.12940>  
542 Ahlström, A. *et al.* (2015). Carbon cycle. The dominant role of semi-arid ecosystems in the trend and  
543 variability of the land CO<sub>2</sub> sink. *Science* **348**, 895–899. <https://doi.org/10.1126/science.aaa1668>

Field Code Changed

544 Araújo, M.B. & New, M. (2007). Ensemble forecasting of species distributions. *Trends Ecol. Evol.* **22**,  
545 42–47. <https://doi.org/10.1016/j.tree.2006.09.010>.

546 Bagstad, K.J. *et al.* (2013). A comparative assessment of decision-support tools for ecosystem services  
547 quantification and valuation. *Ecosyst. Serv.* **5**, 27–39. <https://doi.org/10.1016/j.ecoser.2013.07.004>

548 Barredo, J.I. *et al.* (2012). *A European map of living forest biomass and carbon stock*. (European  
549 Commission, Joint Research Centre). [https://op.europa.eu/en/publication-detail/-](https://op.europa.eu/en/publication-detail/-/publication/b9345574-a96f-4417-87ed-1a85d2252834/language-en)  
550 [/publication/b9345574-a96f-4417-87ed-1a85d2252834/language-en](https://op.europa.eu/en/publication-detail/-/publication/b9345574-a96f-4417-87ed-1a85d2252834/language-en)

551 Bell, V.A. *et al.* (2018a). The MaRIUS- G2G datasets: Grid- to- Grid model estimates of flow and  
552 soil moisture for Great Britain using observed and climate model driving data. *Geosci. Data J.* **5**,  
553 63–72. <https://doi.org/10.1002/gdj3.55>

554 Bell, V.A. *et al.* (2018b). *Grid-to-Grid model estimates of monthly mean flow and soil moisture for*  
555 *Great Britain (1891 to 2015): observed driving data [MaRIUS-G2G-Oudin-monthly]*. [Data Set]  
556 (NERC Environmental Information Data Centre). [https://doi.org/10.5285/f52f012d-9f2e-42cc-](https://doi.org/10.5285/f52f012d-9f2e-42cc-b628-9cdea4fa3ba0)  
557 [b628-9cdea4fa3ba0](https://doi.org/10.5285/f52f012d-9f2e-42cc-b628-9cdea4fa3ba0)

558 Brun, P. *et al.* (2020). Model complexity affects species distribution projections under climate  
559 change. *J. Biogeogr.* **47**, 130–142. <https://doi.org/10.1111/jbi.13734>

560 Bryant, B.P. *et al.* (2018). Transparent and feasible uncertainty assessment adds value to applied  
561 ecosystem services modeling. *Ecosyst.Serv.* **33**, 103–109.  
562 <https://doi.org/10.1016/j.ecoser.2018.09.001>

563 Chaplin-Kramer, R. *et al.* (2019). Global modeling of nature’s contributions to people. *Science* **366**,  
564 255–258. <https://science.sciencemag.org/content/366/6462/255.abstract>

565 Costanza, R. *et al.* (2014). Changes in the global value of ecosystem services. *Glob. Environ.*  
566 *Change* **26**, 152–158. <https://doi.org/10.1016/j.gloenvcha.2014.04.002>

567 Costanza, R. *et al.* (2017). Twenty years of ecosystem services: how far have we come and how far do  
568 we still need to go? *Ecosyst. Serv.* **28**, 1–16. <https://doi.org/10.1016/j.ecoser.2017.09.008>

569 Coxon, G. *et al.* (2019a). DECIPHeR v1: Dynamic fluxEs and Connectivity for Predictions of  
570 Hydrology. *Geosci. Model Dev.* **12**, 2285–2306. <https://doi.org/10.5194/gmd-12-2285-2019>

571 Coxon, G. *et al.* (2019b). *DECIPHeR model estimates of daily flow for 1366 gauged catchments in*  
572 *Great Britain (1962-2015) using observed driving data*. [Data Set] (NERC Environmental  
573 Information Data Centre). <https://doi.org/10.5285/d770b12a-3824-4e40-8da1-930cf9470858>

574 Crossman, N.D., Bryan, B.A. & Summers, D.M. (2012). Identifying priority areas for reducing species  
575 vulnerability to climate change. *Divers. Distrib.* **18**, 60–72. [https://doi.org/10.1111/j.1472-](https://doi.org/10.1111/j.1472-4642.2011.00851.x)  
576 [4642.2011.00851.x](https://doi.org/10.1111/j.1472-4642.2011.00851.x)

577 Diengdoh, V.L. *et al.* (2020). A validated ensemble method for multinomial land-cover  
578 classification. *Ecol. Inform.* **56**, 101065. <https://doi.org/10.1016/j.ecoinf.2020.101065>

579 Ding, H. & Bullock, J.M. (2018). *A Guide to Selecting Ecosystem Service Models for Decision-*  
580 *Making: Lessons from Sub-Saharan Africa*. (World Resources Institute). [wri.org/publication/guide-](http://wri.org/publication/guide-selecting-ecosystem-service)  
581 [selecting-ecosystem-service](http://wri.org/publication/guide-selecting-ecosystem-service)

582 Dormann, C.F. *et al.* (2018). Model averaging in ecology: a review of Bayesian, information-theoretic,  
583 and tactical approaches for predictive inference. *Ecol. Monogr.* **88**, 485–504.  
584 <https://doi.org/10.1002/ecm.1309>

585 Eigenbrod, F. *et al.* (2010) The impact of proxy- based methods on mapping the distribution of  
586 ecosystem services. *J. Appl. Ecol.* **47.2**, 377-385.

587 Elith, J. *et al.* (2011). A statistical explanation of MaxEnt for ecologists. *Divers. Distrib.* **17**, 43–57.  
588 <https://doi.org/10.1111/j.1472-4642.2010.00725.x>

589 Englund, O., Berndes, G. & Cederberg, C. (2017). How to analyse ecosystem services in landscapes—A  
590 systematic review. *Ecol. Indic.* **73**, 492–504. <https://doi.org/10.1016/j.ecolind.2016.10.009>

591 Erceg-Hurn, D.M. & Mirosevich, V.M. (2008). Modern robust statistical methods: an easy way to  
592 maximize the accuracy and power of your research. *Am. Psychol.* **63**, 591–601.  
593 <http://dx.doi.org/10.1037/0003-066X.63.7.591>

Field Code Changed

594 Forestry Commission, United Kingdom. (2018). *National Forest Inventory Woodland GB 2018*. [Data  
595 Set] (Forestry Commission Open Data). [http://data-](http://data-forestry.opendata.arcgis.com/datasets/d3d7bfba1cba4a3b83a948f33c5777c0_0)  
596 [forestry.opendata.arcgis.com/datasets/d3d7bfba1cba4a3b83a948f33c5777c0\\_0](http://data-forestry.opendata.arcgis.com/datasets/d3d7bfba1cba4a3b83a948f33c5777c0_0)  
597 Gassert, F. *et al.* (2015). *Aqueduct Global Maps 2.1*. [Data Set] (World Resources Institute).  
598 <https://www.wri.org/resources/data-sets/aqueduct-global-maps-21-data>  
599 Gobeyn, S. *et al.* (2019). Evolutionary algorithms for species distribution modelling: A review in the  
600 context of machine learning. *Ecol. Modell.* **392**, 179–195.  
601 <https://doi.org/10.1016/j.ecolmodel.2018.11.013>  
602 Graham, C.H. *et al.* (2008). The influence of spatial errors in species occurrence data used in  
603 distribution models. *J Appl. Ecol.* **45**, 239–247. <https://doi.org/10.1111/j.1365-2664.2007.01408.x>  
604 Grenouillet, G. *et al.* (2011). Ensemble modelling of species distribution: the effects of geographical  
605 and environmental ranges. *Ecography* **34**, 9–17. <https://doi.org/10.1111/j.1600-0587.2010.06152.x>  
606 Griggs, D. *et al.* (2013). Sustainable development goals for people and planet. *Nature* **495**, 305–307.  
607 <https://doi.org/10.1038/495305a>  
608 de Groot, R. *et al.* (2012). Global estimates of the value of ecosystems and their services in monetary  
609 units. *Ecosyst. Serv.* **1**, 50–61. <https://doi.org/10.1016/j.ecoser.2012.07.005>  
610 Hao, T. *et al.* (2020). Testing whether ensemble modelling is advantageous for maximising predictive  
611 performance of species distribution models. *Ecography* **43**, 549–558.  
612 <https://doi.org/10.1111/ecog.04890>  
613 He, X. *et al.* (2021). Climate-informed hydrologic modeling and policy typology to guide managed  
614 aquifer recharge. *Science Advances* **7**, p.eabe6025. <https://doi.org/10.1126/sciadv.abe6025>  
615 He, K., Yu, L. & Lai, K.K. (2012). Crude oil price analysis and forecasting using wavelet decomposed  
616 ensemble model. *Energy* **46**, 564–574. <https://doi.org/10.1016/j.energy.2012.07.055>  
617 Henrys, P.A., Keith, A. & Wood, C.M. (2016). *Model estimates of aboveground carbon for Great*  
618 *Britain*. [Data Set] (NERC Environmental Information Data  
619 Centre). <https://doi.org/10.5285/9be652e7-d5ce-44c1-a5fc-8349f76f5f5c>  
620 Inoue, H. & Narihisa, H. (2000) in *Knowledge Discovery and Data Mining. Current Issues and New*  
621 *Applications* (eds Terano, T, Liu, H. & Chen, A.L.P.) 177-180 (Springer).  
622 <https://link.springer.com/book/10.1007/3-540-45571-X>  
623 Kareiva, P. *et al.* (2011). *Natural Capital: Theory and Practice of Mapping Ecosystem Services*.  
624 (Oxford University Press).  
625 <https://oxford.universitypressscholarship.com/view/10.1093/acprof:oso/9780199588992.001.0001/a>  
626 [cprof-9780199588992](https://oxford.universitypressscholarship.com/view/10.1093/acprof:oso/9780199588992)  
627 Keselman, H. J. *et al.* (2008). A generally robust approach for testing hypotheses and setting  
628 confidence intervals for effect sizes. *Psychol. Methods* **13**, 110–129.  
629 <https://doi.apa.org/doi/10.1037/1082-989X.13.2.110>  
630 Kindermann, G.E. *et al.* (2008). A global forest growing stock, biomass and carbon map based on FAO  
631 statistics. *Silva Fennica* **42**, 397–396. <http://pure.iiasa.ac.at/id/eprint/8616/>  
632 Knutti, R., Masson, D. & Gettelman, A. (2013). Climate model genealogy: Generation CMIP5 and  
633 how we got there. *Geophys. Res. Lett.* **40**, 1194–1199. <https://doi.org/10.1002/grl.50256>  
634 Lewis, K.A. *et al.* (2021). Using multiple ecological models to inform environmental decision-  
635 making. *Front. Mar. Sci.* **8**, 283. <https://doi.org/10.3389/fmars.2021.625790>  
636 Liu, D., Li, T. & Liang, D. (2020). An integrated approach towards modeling ranked weights. *Comput.*  
637 *Ind. Eng.* **147**, 106629. <https://doi.org/10.1016/j.cie.2020.106629>  
638 Malinga, R. *et al.* (2015). Mapping ecosystem services across scales and continents—A review. *Ecosyst.*  
639 *Serv.* **13**, 57–63. <https://doi.org/10.1016/j.ecoser.2015.01.006>  
640 Marmion, M. *et al.* (2009). Evaluation of consensus methods in predictive species distribution  
641 modelling. *Divers. Distrib.* **15**, 59–69. <https://doi.org/10.1111/j.1472-4642.2008.00491.x>  
642 Martínez-Harms, M.J. & Balvanera, P. (2012). Methods for mapping ecosystem service supply: a  
643 review. *Int. J. Biodivers. Sci. Ecosyst. Serv. Manag.* **8**, 17-25.

Field Code Changed

644 Martínez-López, J. *et al.* (2019). Towards globally customizable ecosystem service models. *Sci. Total*  
645 *Environ.* **650**, 2325–2336. <https://doi.org/10.1016/j.scitotenv.2018.09.371>

646 Masson, D. & Knutti, R. (2011). Climate model genealogy. *Geophys. Res. Lett.* **38**, L08703.  
647 <https://doi.org/10.1029/2011GL046864>

648 Mulligan M. (2013). WaterWorld: a self-parameterising, physically based model for application in  
649 data-poor but problem-rich environments globally. *Hydrol. Res.* **44**, 748–69.  
650 <https://doi.org/10.2166/nh.2012.217>

651 Ochoa, V. & Urbina-Cardona, N. (2017). Tools for spatially modeling ecosystem services: Publication  
652 trends, conceptual reflections and future challenges. *Ecosyst.Serv.* **26**, 155–169.  
653 <https://doi.org/10.1016/j.ecoser.2017.06.011>

654 Pascual, U. *et al.* (2017). Valuing nature’s contributions to people: the IPBES approach. *Curr. Opin.*  
655 *Environ. Sustain.* **26–27**, 7–16. <https://doi.org/10.1016/j.cosust.2016.12.006>

656 Redhead, J.W. *et al.* (2016). Empirical validation of the InVEST water yield ecosystem service model  
657 at a national scale. *Sci. Total Environ.* **569**, 1418–1426 (2016).  
658 <https://doi.org/10.1016/j.scitotenv.2016.06.227>

659 Refsgaard, J.C. *et al.* (2014). A framework for testing the ability of models to project climate change  
660 and its impacts. *Clim. Change* **122**, 271–282. <https://doi.org/10.1007/s10584-013-0990-2>

661 Scholes, R.J. (1998). *The South African 1: 250 000 maps of areas of homogeneous grazing potential.*  
662 (CSIR, South Africa). No internet reference

663 Sharps, K. *et al.* (2017). Comparing strengths and weaknesses of three ecosystem services modelling  
664 tools in a diverse UK river catchment. *Sci. Total Environ.* **584**, 118–130.  
665 <https://doi.org/10.1016/j.scitotenv.2016.12.160>

666 Smith, B. *et al.* (2014). Implications of incorporating N cycling and N limitations on primary  
667 production in an individual-based dynamic vegetation model. *Biogeosciences* **11**, 2027–2054.  
668 <https://doi.org/10.5194/bg-11-2027-2014>

669 van Soesbergen, A. & Mulligan, M. (2018). Uncertainty in data for hydrological ecosystem services  
670 modelling: Potential implications for estimating services and beneficiaries for the CAZ Madagascar.  
671 *Ecosyst. Serv.* **33**, 175–186. <https://doi.org/10.1016/j.ecoser.2018.08.005>

672 Suich, H., Howe, C. & Mace, G. (2015). Ecosystem services and poverty alleviation: A review of the  
673 empirical links. *Ecosyst. Serv.* **12**, 137–147. <https://doi.org/10.1016/j.ecoser.2015.02.005>

674 Tebaldi, C. & Knutti, R. (2007). The use of the multi-model ensemble in probabilistic climate  
675 projections. *Philos. Trans. R. Soc. A Math. Phys. Eng. Sci.* **365**, 2053–2075.  
676 <https://doi.org/10.1098/rsta.2007.2076>

677 Thibaut, L.M. & Connolly, S.R. (2013). Understanding diversity–stability relationships: towards a  
678 unified model of portfolio effects. *Ecol. Lett.* **16**, 140–150. <https://doi.org/10.1111/ele.12019>

679 Thomas, A. *et al.* (2020). Fragmentation and thresholds in hydrological flow- based ecosystem  
680 services. *Ecol. Appl.* **30**, e02046. <https://doi.org/10.1002/eap.2046>

681 UKNEA. (2011). *The UK National Ecosystem Assessment: Synthesis of the Key Findings.* (UNEP-  
682 WCMC, Cambridge). [https://www.unep-wcmc.org/resources-and-data/UK-national-ecosystem-](https://www.unep-wcmc.org/resources-and-data/UK-national-ecosystem-assessment)  
683 [assessment](https://www.unep-wcmc.org/resources-and-data/UK-national-ecosystem-assessment)

684 Verhagen, W. *et al.* (2017). Use of demand for and spatial flow of ecosystem services to identify  
685 priority areas. *Conserv. Biol.* **31**, 860–871. <https://doi.org/10.1111/cobi.12872>

686 Wang, H.M. *et al.* (2019). Does the weighting of climate simulations result in a better quantification of  
687 hydrological impacts? *Hydrol. Earth Syst. Sci.* **23**, 4033–4050. [https://doi.org/10.5194/hess-23-](https://doi.org/10.5194/hess-23-4033-2019)  
688 [4033-2019](https://doi.org/10.5194/hess-23-4033-2019)

689 Willcock, S. *et al.* (2016). Do ecosystem service maps and models meet stakeholders’ needs? A  
690 preliminary survey across sub-Saharan Africa. *Ecosyst. Serv.* **18**, 110–117.  
691 <https://doi.org/10.1016/j.ecoser.2016.02.038>

692 Willcock, S. *et al.* (2019). A Continental-Scale Validation of Ecosystem Service  
693 Models. *Ecosystems* **22**, 1902–1917. <https://doi.org/10.1007/s10021-019-00380-y>

Field Code Changed

694 Willcock, S. *et al.* (2020). Ensembles of ecosystem service models can improve accuracy and indicate  
695 uncertainty. *Sci. Total Environ.* **747**, 141006. <https://doi.org/10.1016/j.scitotenv.2020.141006>  
696 Wong, C.P. *et al.* (2014). Linking ecosystem characteristics to final ecosystem services for public  
697 policy. *Ecol. Lett.* **18**, 108–118. <https://doi.org/10.1111/ele.12389>  
698 Zulian, G. *et al.* (2018). Practical application of spatial ecosystem service models to aid decision  
699 support. *Ecosyst. Serv.* **29**, 465–480. <https://doi.org/10.1016/j.ecoser.2017.11.005>  
700

Field Code Changed

Field Code Changed

## Response to Comments from the Editors and Reviewers

### **Comments Editor-in-Chief:**

You may consider adjusting the title to make it more easily understandable to non-modellers.  
for example: "Reducing Uncertainty in Ecosystem Service Modelling through Weighted Ensembles"

RESPONSE: We have made this change

### **Comments Associate editor:**

I am pleased to accept this highly relevant manuscript for publication pending some very minor changes. First, please avoid using abbreviations in your highlights. Second, please go over the suggestions made by Reviewer 1 and correct the language where needed.

RESPONSE: The highlights have been amended as requested. We have address all of R1's comments (below). We thank both the editors and reviewers for their helpful suggestions during the review process.

### **Comments Reviewer #1:**

The authors have addressed all my comments and questions in a thoughtful and courteous way, and I think this paper is in good shape for publication. I still think it makes a significant contribution to the literature in terms of providing evidence for the accuracy of ecosystem service models using validation data, and advances methods for ensemble ES modeling.

RESPONSE: Thank you

I have a few very minor edits to recommend:

Lines 89-90

Should be framed as a question with quotation marks:

Thus, here we explore the outstanding question, "what are the best ways... sustainability science?"

RESPONSE: This has been done

Lines 104-106

This sentence is awkwardly worded, I recommend rewording e.g. moving the parenthetical to the end of the sentence

RESPONSE: This sentence has been changed from:

*"As well as varying considerably in their underlying method, ES models often differ in the forms of their outputs (e.g. summed monetary value of the ES (de Groot et al. 2012) vs. specific biophysical predictions), even when modelling the same ES"*

To:

*"As well as varying considerably in their underlying method, ES models often differ in the forms of their outputs, even when modelling the same ES (e.g. summed monetary value of the ES (de Groot et al. 2012) vs. specific biophysical predictions)."*

Lines 107-112

There are two run-on sentences here (starting "An important knowledge gap... and "Since models for a particular ES...") - I recommend splitting them up

RESPONSE: These sentences have been split up. This section now reads:

*"An important knowledge gap is therefore how to combine distinct ES model outputs as complementary inputs to provide a reliable ensemble. Outputs from different ES models can have different units and it is challenging to decide the relative weighting to place on each model. Models for a particular ES often have different structures, may include different processes, or may represent the same processes in different ways (Ochoa & Urbina-Cardona 2017). As a result, the different ES models will most likely not have equal accuracy, and so prediction errors (i.e. bias) may not be normally distributed among models (Dormann et al. 2018)."*

159

"ArcPy" should be capitalized, I believe

RESPONSE: This has been done throughout

184 (and multiple places)

I have never come across "gridcell" used as a single word, unless this is common I recommend two words (throughout the paper)

RESPONSE: This has been changed to two words throughout

Table 2

"Grid size (grain)"

I have more commonly seen this referred to as "resolution" or "spatial resolution", I would use that term here or somewhere in the text, I have not heard the term "grain" used this way before.

RESPONSE: We have replaced 'grain' with 'spatial resolution' throughout.

342

"However, we are aware of introducing a potential bias that could skew non-forested areas to lower values."

Please explain, as currently this sentence seems rather out of place.

RESPONSE: We have added further explanation to this sentence:

*"However, since our validation data are only from forests/woodlands, we are aware of introducing a potential bias that could skew non-forested areas to lower values."*

425

"models have substantial higher" should read "substantially higher"

RESPONSE: Change made

452

""However, this effect is lessened if models share similar assumptions and, therefore, concomitant biases - highlighting the importance of including multiple model outputs"

- I would think that similar biases/assumptions between models highlights the importance of using validation data, not the importance of using multiple model outputs?

RESPONSE: We have added this to the sentence:

*"However, this effect is lessened if models share similar assumptions and, therefore, concomitant biases - highlighting the importance of including multiple model outputs and, where data are available, model validation."*

475

"likely because of over-fitting to woodland areas"

- Can you provide just a little more explanation of this, why is the over-fitting to woodland areas and not other habitat types? Is this the reason why there's a potential bias that could skew non-forested areas to lower values (line 342, above)?

RESPONSE: We have added further explanation as follows:

*"As a result accuracy-derived weights show high uncertainty in areas where training data were not available (i.e. non-forested areas; SI-4), likely because of over-fitting to areas with available data (i.e. forests/woodlands) producing correlative patterns that explain other areas less well."*

485-86

"Our results should serve as a 'call to arms' for ES researchers and practitioners to increasingly use ensembles of models to support decision-making for sustainability."

I would like to see a similar statement made in the abstract, as a key take-away of the paper.

RESPONSE: Word limits prevent us adding more to the abstract. But we have strengthened the last sentence in order to convey a similar message:

*"To support robust decision-making for sustainable development and reducing uncertainty around these decisions, our analysis suggests various ensemble methods should be applied depending on data quality, for example if validation data are available."*

502

How many models, at a minimum, would the authors count as an "ensemble"? Two? More than two?

RESPONSE: Anything >1 could be considered an ensemble. But the optimum number of models to include in the ensemble will be context specific, and we would not be comfortable speculating on that here since such would require a marginal gain analysis, which is beyond the scope of this manuscript .

498

"The most complex (a priori best) ES models require substantial inputs (i.e. data, computational power, subscription fees, and staff time), and so running multiple models - whilst requiring additional resources - results in a large gain per extra unit resource."

- I understand the argument (that complex ES models already require substantial time, so why not run additional models?) - but I still don't completely buy it. I laud the authors for making their data freely available, as this itself will be the biggest contribution for those who don't have the time, data, or capacity to run ES ensembles themselves. And I am convinced by their premise that ensemble models out-perform individual models, on average. There just aren't multiple models (nor validation data) available for most ES beyond carbon and water (as the authors now acknowledge). And the requirements of running multiple models are substantial, and virtually never feasible outside of academic research. I am not requesting further changes, just pushing back on the practicality of this suggestion for most applications outside of academia. I agree with the author's last statement in this paragraph (multiple models be developed for ES where they are lacking, those who can should share their data freely, etc.)

RESPONSE: We are all in agreement here. As the reviewer suggests, many outside academia will struggle to run ES ensembles (even though they convincingly out-perform individual models). The solutions to this are provided in our last statement, as the reviewer acknowledges:

*"Thus, we recommend that multiple models be developed for ES where they are lacking (e.g. cultural services; Martínez-Harms and Balvanera, 2012; Wong et al. 2014), and that those with access to sufficient resources to run multiple models ensure the ensemble outputs are freely available, making the use of these ensembles more feasible and accessible for all (Willcock et al. 2020)."*

We are already working on addressing this issue, using the established techniques in this manuscript to create ES ensembles at a global-scale for carbon, water, sediment retention, recreation, grazing and fuelwood. Once complete, we will make these layers publicly available, and so support the use of ensembles outside academia. This is work in progress, but out of the scope of this manuscript.

**Reviewer #2:**

I appreciate the authors' efforts in addressing my comments. Thank you!

I recommend the paper to be published.

RESPONSE: Thank you.

## **TITLE PAGE**

# **Reducing Uncertainty in Ecosystem Service Modelling through Weighted Ensembles**

Hooftman, Danny A.P.<sup>1,2</sup>, James M. Bullock<sup>2</sup>, Laurence Jones<sup>3</sup>, Felix Eigenbrod<sup>4</sup>, José I. Barredo<sup>5</sup>, Matthew Forrest<sup>6</sup>, Georg Kindermann<sup>7</sup>, Amy Thomas<sup>3</sup> & Simon Willcock<sup>8,9\*</sup>

\* Corresponding author

### **Affiliations:**

1. Lactuca: Environmental Data Analyses and Modelling, The Netherlands. [danny.hooftman@lactuca.nl](mailto:danny.hooftman@lactuca.nl)
2. UK Centre for Ecology and Hydrology, Wallingford, OX10 8BB, United Kingdom. [jmbul@ceh.ac.uk](mailto:jmbul@ceh.ac.uk)
3. UK Centre for Ecology and Hydrology, Bangor, LL57 2UW, United Kingdom. [lj@ceh.ac.uk](mailto:lj@ceh.ac.uk), [athomas@ceh.ac.uk](mailto:athomas@ceh.ac.uk)
4. Geography and Environment, University of Southampton, United Kingdom. [F.Eigenbrod@soton.ac.uk](mailto:F.Eigenbrod@soton.ac.uk)
5. Joint Research Centre of the European Commission, Brussels, Belgium. [Jose.BARREDO@ec.europa.eu](mailto:Jose.BARREDO@ec.europa.eu)
6. Senckenberg Biodiversity and Climate Research Centre, Frankfurt, Germany. [matthew.forrest@senckenberg.de](mailto:matthew.forrest@senckenberg.de)
7. International Institute for Applied Systems Analysis, Laxenburg, Austria. [kinder@iiasa.ac.at](mailto:kinder@iiasa.ac.at)
8. School of Natural Sciences, Bangor University, United Kingdom. [s.willcock@bangor.ac.uk](mailto:s.willcock@bangor.ac.uk)
9. Rothamsted Research, Harpenden, Hertfordshire, AL5 2JQ, United Kingdom.

**Contributions:** DAPH, JMB & SW conceived the project. DAPH, LJ, AT, MF, JB & GK provided ES model descriptions and outputs. DAPH conducted all analyses. DAPH, JMB & SW wrote the manuscript, with comments from AT, FE, JB, LJ, MF & GK.

**Acknowledgements:** This work took place under the Ensembles project – Using ensemble techniques to capture the accuracy and sensitivity of ecosystem service models ([NE/T00391X/1](#)). Land Cover Map 2015 is under UKCEH licence 1403. We acknowledge the help of Kevin Watts for guiding us through the Forest Research data and John Redhead for providing InVEST biophysical tables. We also thank the anonymous reviewers for their insightful comments on the manuscript.

## **TITLE PAGE**

# **Weighted Ensembles Reduce Uncertainty in Ecosystem Service Modelling** **Reducing Uncertainty in Ecosystem Service Modelling through Weighted Ensembles**

Hoofman, Danny A.P.<sup>1,2</sup>, James M. Bullock<sup>2</sup>, Laurence Jones<sup>3</sup>, Felix Eigenbrod<sup>4</sup>, José I. Barredo<sup>5</sup>, Matthew Forrest<sup>6</sup>, Georg Kindermann<sup>7</sup>, Amy Thomas<sup>3</sup> & Simon Willcock<sup>8,9\*</sup>

\* Corresponding author

### **Affiliations:**

1. Lactuca: Environmental Data Analyses and Modelling, The Netherlands. [danny.hoofman@lactuca.nl](mailto:danny.hoofman@lactuca.nl)
2. UK Centre for Ecology and Hydrology, Wallingford, OX10 8BB, United Kingdom. [jmbul@ceh.ac.uk](mailto:jmbul@ceh.ac.uk)
3. UK Centre for Ecology and Hydrology, Bangor, LL57 2UW, United Kingdom. [lj@ceh.ac.uk](mailto:lj@ceh.ac.uk), [athomas@ceh.ac.uk](mailto:athomas@ceh.ac.uk)
4. Geography and Environment, University of Southampton, United Kingdom. [F.Eigenbrod@soton.ac.uk](mailto:F.Eigenbrod@soton.ac.uk)
5. Joint Research Centre of the European Commission, Brussels, Belgium. [Jose.BARREDO@ec.europa.eu](mailto:Jose.BARREDO@ec.europa.eu)
6. Senckenberg Biodiversity and Climate Research Centre, Frankfurt, Germany. [matthew.forrest@senckenberg.de](mailto:matthew.forrest@senckenberg.de)
7. International Institute for Applied Systems Analysis, Laxenburg, Austria. [kinder@iiasa.ac.at](mailto:kinder@iiasa.ac.at)
8. School of Natural Sciences, Bangor University, United Kingdom. [s.willcock@bangor.ac.uk](mailto:s.willcock@bangor.ac.uk)
9. Rothamsted Research, Harpenden, Hertfordshire, AL5 2JQ, United Kingdom.

**Contributions:** DAPH, JMB & SW conceived the project. DAPH, LJ, AT, MF, JB & GK provided ES model descriptions and outputs. DAPH conducted all analyses. DAPH, JMB & SW wrote the manuscript, with comments from AT, FE, JB, LJ, MF & GK.

**Acknowledgements:** This work took place under the Ensembles project – Using ensemble techniques to capture the accuracy and sensitivity of ecosystem service models ([NE/T00391X/1](#)). Land Cover Map 2015 is under UKCEH licence 1403. We acknowledge the help of Kevin Watts for guiding us through the Forest Research data and John Redhead for providing InVEST biophysical tables. We also thank the anonymous reviewers for their insightful comments on the manuscript.



[Click here to access/download](#)

**Video**

Ensembles\_of\_ecosystem\_service\_models - Final no  
subs.mp4

