

1 **Mapping habitat indices across river networks using spatial statistical modelling of**
2 **River Habitat Survey data.**

3 Marc Naura¹, Mike J. Clark² (†), David A. Sear², Peter M. Atkinson³, Duncan D. Hornby²,
4 Paul Kemp¹, Judy England⁴, Graeme Peirson⁵, Chris Bromley⁶, Matthew G. Carter⁷

5 ¹ Faculty of Engineering and the Environment, University of Southampton, SO17 1BJ,
6 Southampton, UK; ² School of Geography, University of Southampton, Southampton,
7 SO17 1BJ, UK; ³ Faculty of Science and Technology, Engineering Building, Lancaster
8 University, Bailrigg, Lancaster LA1 4YR, UK; ⁴ Environment Agency, Red Kite House,
9 Howbery Park, Crowmarsh Gifford, Wallingford OX10 8BD, UK; ⁵ Mance House, Arthur
10 Drive, Hoo Farm Industrial Estate, Worcester Road, Kidderminster DY11 7RA, UK. ⁶
11 Scottish Environment Protection Agency, Strathallan House, Castle Business Park,
12 Stirling FK9 4TZ, UK; ⁷ Environment Agency, Thames East, Apollo Court, 2 Bishops
13 Square, St. Albans, Road West, Hatfield, Hertfordshire AL10 9EX, UK.

14 Corresponding author: Marc Naura, marc.naura@soton.ac.uk, Faculty of Engineering and
15 the Environment, University of Southampton, SO17 1BJ, Southampton; tel: + 44 2380
16 592700

17 **Abstract**

18 Freshwater ecosystems are declining faster than their terrestrial and marine
19 counterparts because of physical pressures on habitats. European legislation requires
20 member states to achieve ecological targets through the effective management of
21 freshwater habitats. Maps of habitats across river networks would help diagnose
22 environmental problems and plan for the delivery of improvement work. Existing
23 habitat mapping methods are generally time consuming, require experts and are
24 expensive to implement. Surveys based on sampling are cheaper but provide patchy

25 representations of habitat distribution. In this study, we present a method for mapping
26 habitat indices across networks using semi-quantitative data and a geostatistical
27 technique called regression kriging. The method consists of the derivation of habitat
28 indices using multivariate statistical techniques that are regressed on map-based
29 covariates such as altitude, slope and geology. Regression kriging combines the
30 Generalised Least Squares (GLS) regression technique with a spatial analysis of
31 model residuals. Predictions from the GLS model are 'corrected' using weighted
32 averages of model residuals following an analysis of spatial correlation. The method
33 was applied to channel substrate data from the River Habitat Survey in Great Britain.
34 A Channel Substrate Index (CSI) was derived using Correspondence Analysis and
35 predicted using regression kriging. The model explained 74% of the main sample
36 variability and 64% in a test sample. The model was applied to the English and Welsh
37 river network and a map of CSI was produced. The proposed approach demonstrates
38 how existing national monitoring data and geostatistical techniques can be used to
39 produce continuous maps of habitat indices at the national scale.

40 **Keywords:** habitat mapping, habitat indices, channel substrate, regression kriging, River
41 Habitat Survey, geostatistics

42

43 **1. Introduction**

44 Freshwater ecosystems represent less than 1% of the Earth's surface and 10% of all
45 known species, yet they are declining faster and are more endangered than their
46 terrestrial or marine counterparts, partly because of physical pressures on habitats and
47 species (Loh et al., 2005; Revenga et al., 2005; Strayer and Dudgeon, 2010; Vorosmarty
48 et al., 2010; WWF, 2014).

49 Although research in ecology and environmental management has grown substantially in
50 the past half-century, it has mainly focused on post-industrial issues such as water
51 quality, pollution and land use impacts (Vaughan et al., 2009). With gradual improvement
52 in water quality, other limiting factors such as physical habitat quality (i.e. the naturalness
53 of the flow of water, and the structure and composition of the river bed and banks) and
54 connectivity have become prominent.

55 Globally, degradation of physical habitat quality due to river engineering and associated
56 activities (e.g. constructions of dams, bridges, concrete banks, dredging) is recognised as
57 a major conservation issue (Collen et al., 2014; Sala et al., 2000; Tockner and Stanford,
58 2002; World Conservation Monitoring Centre, 1998). In Europe, as part of the
59 implementation of the Water Framework Directive (WFD), member states must assess
60 the ecological condition of rivers and lakes based on the naturalness of a series of
61 biological elements (European Union, 2000). Following the first round of River Basin
62 Management Planning, 56% of water bodies failed to achieve their ecological targets.
63 Engineered structures and 'altered habitats' were the dominant pressures responsible for
64 the failure, ahead of point and diffuse sources of pollution (European Environment
65 Agency, 2012). In England and Scotland, the proportion of water bodies failing to achieve
66 ecological targets because of physical alterations was 49% and 37%, respectively
67 (Environment Agency, 2012). The WFD requires member states to mitigate or remove

68 impacts on habitats and species through the implementation of programmes of measures
69 including river restoration.

70 The effective management of habitats at global and local scales should ideally be based
71 on some knowledge of their distribution and an assessment of their naturalness and
72 accessibility. At present, in Great Britain, habitats are either surveyed using semi-
73 quantitative methods at randomly selected sites that do not allow for continuous
74 assessments or using habitat mapping techniques over longer stretches of river
75 (Maddock, 1999). Habitat mapping is geographically limited and generally carried out on
76 an *ad hoc* basis by experts during 'walkover surveys' where habitat features are recorded
77 on maps using mobile Geographic Information System (GIS) or hand-drawn sketches and
78 some broad typologies (Hendry and Cragg-Hine, 1997; Sear et al., 2009). Although such
79 methods provide valuable information on habitat distributions over relatively small areas,
80 they are likely to be too expensive to implement across entire networks. The reliance on
81 expert judgement for assessing habitat types and boundaries may also generate
82 between-surveyor variability in the outputs produced and, as notions of habitat structure
83 evolve, data collected at one point in time may not be comparable to maps produced
84 years later by different experts (Cherrill and McClean, 1999).

85 An alternative approach is to use river typologies based on geomorphological templates
86 to predict the occurrence of broad river types along the river continuum. The history of
87 attempts to classify rivers into different types spans at least 125 years, a period over
88 which perhaps a hundred if not more individual efforts to divide and categorise rivers have
89 been made (reviews of the extent of such efforts are given by Downs, 1995; Montgomery
90 and Buffington, 1997; Mosley, 1987; Naiman et al., 1992; Newson et al., 1998; Thorne,
91 1997).

92 Most river classification systems are based on the identification of river types using a few
93 key variables representing drivers of geomorphological change or river processes such as
94 stream power, sediment transport and supply (Montgomery and Buffington, 1997;
95 Newson et al., 1998; Rosgen, 1994). Although relationships between expert-driven
96 geomorphic types and GIS attributes such as slope and drainage area can be observed,
97 there is a considerable amount of overlap between types, reflecting the potential influence
98 of additional driving elements such as channel, bank and hillslope vegetation, climate,
99 woody debris, and natural variability in channel process expression (Church, 2002;
100 Montgomery and Buffington, 1997; Rosgen, 1994). Greater differentiation between river
101 types can be achieved by introducing attributes recorded in the field such as relative
102 roughness (Montgomery and Buffington, 1997), shear stress or channel substrate
103 (Rosgen, 1994), but this implies that extensive field work is carried out, thus reducing the
104 feasibility of such an approach at national scales.

105 In this article, we propose an alternative approach for mapping habitat elements across
106 entire river networks that does not require continuous surveys of river catchments, but
107 makes use of existing semi-quantitative survey data, GIS and a geostatistical technique
108 called regression kriging (RK). The principle of the method is to identify and define habitat
109 indices representing major dimensions in habitat distribution using known equations,
110 expert systems or multivariate statistical analysis applied to existing habitat data taken
111 from national surveys or monitoring programmes. The habitat indices are then predicted
112 using Generalised Least Squares (GLS) linear regression models using GIS map-derived
113 covariates such as altitude, slope, distance from source, discharge and geology which
114 represent the known drivers of habitat/geomorphological change. The model residuals are
115 then analysed using geostatistical functions to identify any remaining spatial correlation
116 and pattern in their distribution. In the presence of spatial correlation, an interpolation

117 method, called kriging, is applied to account for (and, thus remove) any spatially
118 correlated residual variance such that the interpolated residual predictions can be added
119 to the GLS regression predictions. The RK model can then be applied to the entire river
120 network by deriving the GIS covariates at regular spatial intervals (e.g. 500 m).

121 This paper reports the development and application of the statistical models to a key and
122 poorly mapped habitat element – channel substrate. Channel substrate is a key
123 component of species habitat (Maddock, 1999; Townsend and Hildrew, 1994), and it is
124 one of three elements defining morphological condition under the WFD (European Union,
125 2000). Channel substrate is also linked to the wider issues of diffuse pollution and
126 agricultural impacts and it is key to our understanding of river and catchment processes
127 (Collins et al., 2014; Rosgen, 1994).

128 **2. Material and methods**

129 2.1. Index derivation

130 River Habitat Survey (RHS) data was used to derive an index representing channel
131 substrate. RHS is a CEN-compliant (CEN, 2004) standard methodology for
132 hydromorphological assessment under the WFD that is used in the UK and across
133 Europe (Raven et al., 1997). It is a methodology for recording habitat features for wildlife
134 that has been implemented at more than 25,000 sites in the UK since 1994. From 1994
135 to 1996 and from 2007 to 2008, surveys were carried out at random sites in every 10 km²
136 in England and Wales, thus, ensuring a wide geographical coverage of the river network.

137 RHS records the presence of natural and management features at 10 equally spaced
138 transects or 'spot-checks' along a 500 m reach (Raven et al., 1997). A visual estimate of
139 the dominant channel surface substrate classified into eight categories according to the
140 Wentworth scale (Wentworth, 1922) is recorded at each spot-check. The substrate types

141 recorded (with acronyms in brackets) are bedrock (BE), boulder (BO), cobble (CO),
142 gravel-pebble (GP), sand (SA), silt (SI), clay (CL) and peat (PE). When channel substrate
143 is not visible because of depth, water turbidity or the presence of a culvert, surveyors
144 record the substrate type as 'Not Visible' (NV).

145 RHS spot-check data on channel substrate was tabulated for all existing sites, each row
146 representing a site and each column a substrate type (including 'Not visible'). The
147 channel substrate spot-check table was analysed using Correspondence Analysis (CA).
148 CA is a multivariate analytical technique similar to Principal Component Analysis that is
149 applicable to contingency tables (i.e. tables of counts). CA performs an analysis of the
150 total table inertia and extracts dimensions (or components) representing linear
151 combinations of input variables based on the amount of total inertia explained. Only sites
152 in Great Britain were used as GIS datasets were not available for Northern Ireland at the
153 time of the analyses.

154 To derive the index, we used a subset of 2680 semi-natural RHS sites (i.e. sites with few
155 or no in-channel bank structures or modifications) to reduce the potential influence of
156 modifications on natural channel substrate diversity (Raven et al., 1997). Missing ('Not
157 Visible') values were added as an additional variable in the analyses to account for
158 differences in survey counts when present. The resulting dimensions were investigated
159 for their ecological and geomorphological significance and for the amount of variability
160 (i.e. inertia) they explained. One dimension was chosen to represent substrate and
161 calculate an index score, called the Channel Substrate Index (CSI) for all sites in the RHS
162 database.

2.2. Regression kriging

RK was applied following an iterative procedure using both Ordinary Least Square (OLS) and GLS regression techniques (Bivand et al., 2008; Webster and Oliver, 2007). The CSI index was first transformed using a Box-Cox procedure and modelled against a series of GIS attributes: four Principal Component Axes (PCA) combining altitude, slope, distance to source and height of source that were shown by Jeffers (1998) and Vaughan et al. (2013) to be strongly correlated to sediment distribution; land use categories from the Land Cover Map 2000 (Fuller et al., 2002); British Geological Survey solid and drift geology categories taken from the 1/625,000 scale maps; hydrometric areas corresponding to large catchment areas; and solid geology age categorised in 11 groups from the pre-Cambrian to the Neogene. Solid geology age was included as a surrogate for hardness as older rocks tend to be harder and display coarser substrate types than softer and younger sedimentary deposits.

Nominal attributes such as solid geology, hydrometric area and land use were transformed into binary indicator variables. Due to the resulting large number of indicator variables which would have rendered the predictive models difficult to display and interpret (e.g. there are more than 100 different solid geology types), indicator variables were grouped based on their relationships to the CSI. Grouping was done by comparing coefficient values of indicator variables when individually regressed against CSI or performing ANOVAs.

Only RHS sites with no missing channel substrate spot-check records were retained for the analysis as their presence introduces a potential bias in channel substrate representation and prediction. Model selection was performed using the Minitab 16 (Minitab, 2010) linear regression (OLS) best subset selection procedure using Mallows Cp (all models) on RHS sites from 1994 to 2005 (9473 sites).

188 Model residuals were analysed for the presence of spatial correlation using a variogram
189 (Webster and Oliver, 2007), which plots semivariance (a measure of dissimilarity) against
190 lag vector (the distance and direction of separation). Spatially uncorrelated data display
191 no observable change in semivariance with an increase in lag distance and are typically
192 represented by a flat variogram. Spatially correlated data are typically represented by a
193 monotonically increasing semivariance as the lag distance between sites increases.

194 The empirical variogram is first calculated as the average squared difference between
195 pairs of data points at each of a series of lags. It is fitted with a permissible variogram
196 model (the model must not result in negative prediction variances) to describe the shape
197 of the curve and identify the parameters which are required for RK. Of particular
198 relevance are the nugget and the range parameters (assuming that a bounded model
199 such as the spherical or exponential model is fitted). The nugget variance is equal to the
200 variance for sites re-surveyed or re-sampled at the same location and expresses micro-
201 scale variability and survey error. The range is the distance at which the semivariance
202 reaches a plateau and beyond which data are no longer spatially correlated (Webster and
203 Oliver, 2007).

204 The OLS variogram parameters were used as part of an iterative process to estimate the
205 regression parameters using GLS. GLS is preferred to OLS as the latter assumes
206 independence of observations, an unlikely case given spatial correlation (Bivand et al.,
207 2008). Model residuals were checked for the presence of trends and outliers. The models
208 were validated using a leave-one-out cross-validation technique (Bivand et al., 2008;
209 Vaughan and Ormerod, 2003).

210 The model was applied to all points using the GSTAT kriging procedure. Kriging linearly
211 averages the residual values of surrounding points with weights estimated using the
212 residual variogram (Webster and Oliver, 2007). To reduce processing time, only points

213 within a radius roughly equal to the distance at which no observable correlation exists (i.e.
214 the range) were selected for kriging. The kriged residuals were then saved and added to
215 the cross-validated predictions from the previous model. A pseudo- R^2 value for the final
216 predictions was derived by correlating the predicted and observed values. Residuals were
217 computed and checked for their distribution and for signs of remaining spatial correlation.
218 A check of model stability against time was performed by examining the residual
219 distribution for each year of survey. The model was tested on a sample of 3884 sites
220 collected from 2006-2011.

221 The predictive model and kriging procedure were applied to the entire English and Welsh
222 1:50,000 river network to produce a national map of sediment distribution. To do so,
223 points were generated every 500 m on the river network using RivEX (Hornby, 2010) and
224 the GIS map-based covariates required for prediction were derived for each point.

225 **3. Results**

226 3.1. Channel Substrate Index

227 The first two components of the CA explained 21% and 17% of the total inertia (Table 1).
228 The first component represented a gradient between sites dominated by fine substrate
229 such as silt, clay and sand, and sites dominated by coarse substrate such as bedrock and
230 boulders (Fig. 1). The first component was defined by the relative occurrence of all
231 substrate types with a greater contribution from silt which explained 35% of the
232 component inertia (Table 1). The first component explained 42% of silt distribution inertia,
233 31% of boulder inertia, 23% for cobbles and 20% for gravel pebble and bedrock.

234 Table 1: Simple CA on channel substrate types for 2680 semi-natural RHS in Great
 235 Britain. Only detailed results for the first 2 components are displayed. The 'Coord'
 236 columns contain the principal coordinates for each substrate type and axis. The 'Contr'
 237 column expresses the relative contribution of individual substrate types to axis definition
 238 whilst the 'Corr' column (or relative contribution) represents the amount of individual
 239 substrate inertia explained by each component (Greenacre, 1993).

240 Individual axes inertia relative to total inertia

Axis	Inertia	Proportion	Cumulative
1	0.7767	0.2051	0.2051
2	0.6442	0.1701	0.3752
3	0.5991	0.1582	0.5334
4	0.5632	0.1487	0.6821
5	0.4317	0.1140	0.7961
6	0.4134	0.1092	0.9052
7	0.3589	0.0948	1.0000

241 Total 3.7872

242
 243 Column Contributions for components 1 and 2

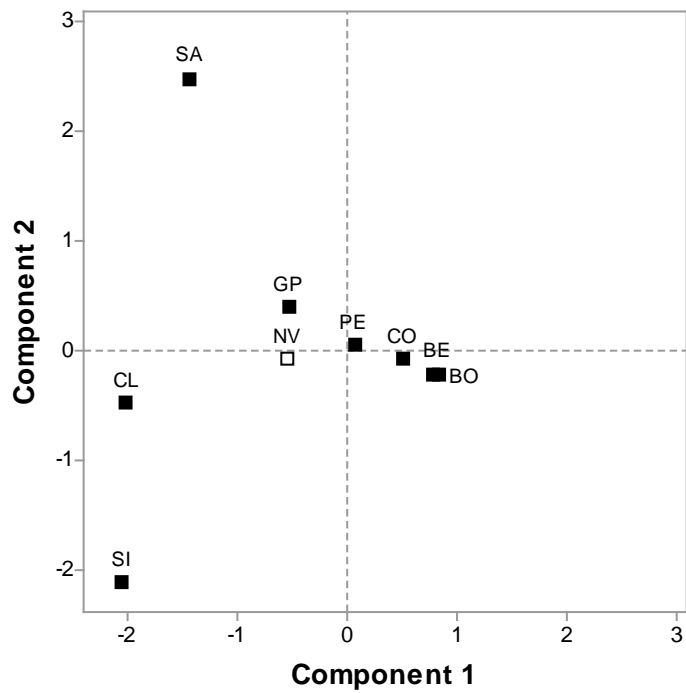
Name	Component .1			Component 2		
	Coord	Corr	Contr	Coord	Corr	Contr
BE	0.793	0.190	0.097	-0.229	0.016	0.010
BO	0.835	0.307	0.151	-0.230	0.023	0.014
CO	0.511	0.226	0.097	-0.077	0.005	0.003
GP	-0.532	0.197	0.106	0.394	0.108	0.070
SA	-1.433	0.160	0.125	2.480	0.480	0.449
SI	-2.052	0.423	0.349	-2.119	0.451	0.449
CL	-2.008	0.097	0.074	-0.483	0.006	0.005
PE	0.067	0.000	0.000	0.039	0.000	0.000

245
 246 Supplementary Columns

Name	Component 1			Component 2		
	Coord	Corr	Contr	Coord	Corr	Contr
NV	-0.540	0.014	0.044	-0.084	0.000	0.001

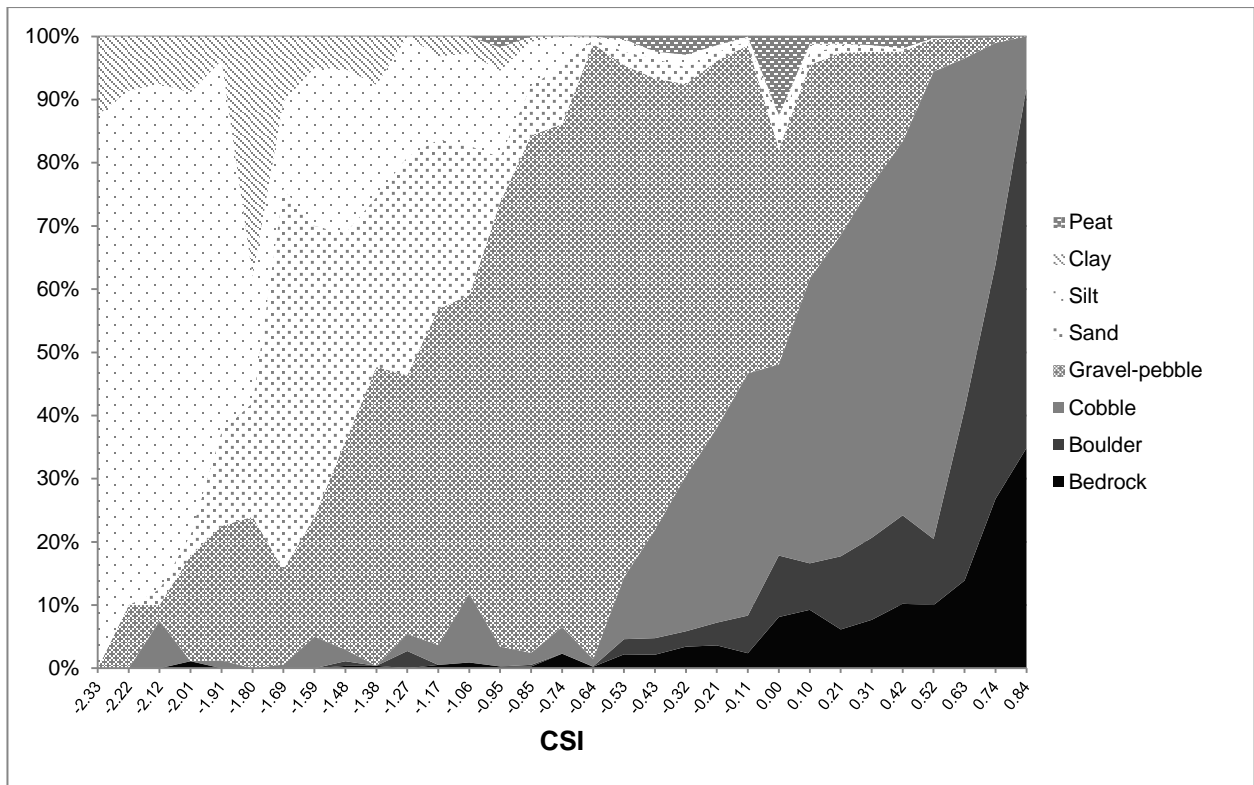
248
 249
 250 The first CA axis represented a gradual increase in substrate size with a gradual shift
 251 from sites dominated by fine sediment to sites dominated by larger substrate (Fig. 2). The
 252 second component represented a gradient between silt and sand dominated sites and
 253 explained nearly 50% of the inertias of both substrate types (Table 1). The remaining
 254 components either represented gradients between two or three substrate types or were

255 linked to the occurrence of rare types such as peat or clay. Missing values were not
256 associated with any particular substrate category and only 1% of the missing values
257 inertia was explained by the first two components.



258

259 Fig. 1: Symmetrical plot of substrate category profiles for the first two CA axes.



260

261 Fig. 2: RHS sites were grouped into 31 bins based on their CSI index values. The plot
 262 displays for each bin the average occurrence of channel substrate types.

263 The first component was chosen for its geomorphological relevance as it represented a
 264 well-known dimension in sediment fining and sorting along the river network (Morris and
 265 Williams, 1999) and has habitat significance with regards to species distribution
 266 (Chessman et al., 2006; Gasparini et al., 1999; Rice et al., 2001). The CSI was calculated
 267 for all existing RHS sites using channel substrate standard coordinates for the first
 268 component in the following equation:

269
$$CSI = (0.89(AR+BE) + 0.95 BO + 0.58 CO + 0.08 PE - 0.6 GP - 1.63 SA - 2.33 SI - 2.28$$

 270
$$CL) / N_{sc}$$

271 where each two-letter acronym refers to RHS channel substrate categories and N_{sc} is the
 272 total number of spot-checks. Artificial channel substrate (AR) was given the same
 273 coefficient as bedrock substrate.

274 3.2. Variable selection

275 Only attributes selected by the best subset procedure will be presented and discussed.
276 Land use categories and drift geologies were not selected in any of the models extracted
277 using the best subset procedure. The attributes retained for the analyses were the PCA
278 axes, solid geology age and categories and hydrometric areas.

279 The four PCA variables represent environmental gradients describing site location and
280 profile (PCA1; i.e. lowland low altitude and slope, and upland high altitude and slope),
281 catchment area (PCA2), local discontinuities in profile/geology (PCA3) and catchment
282 slope (PCA4) (Jeffers, 1998).

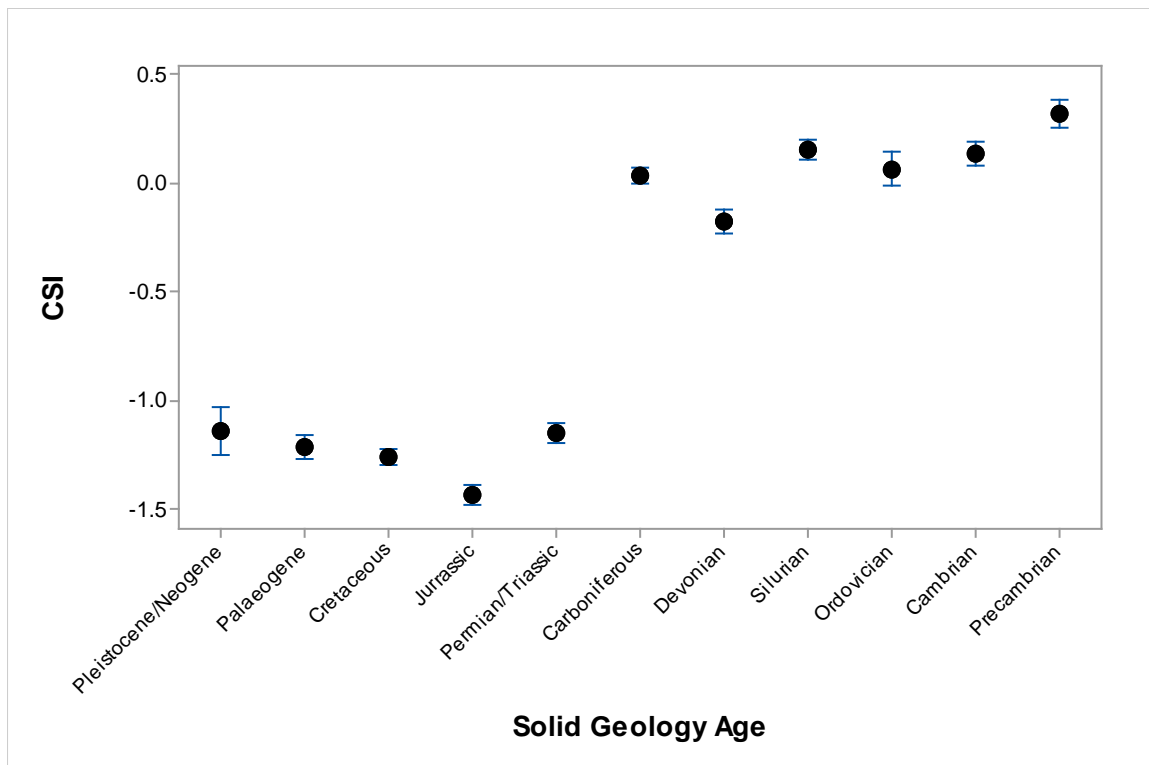
283 Solid geologies were split into two groups based on their age. More recent geologies from
284 the Permian and Triassic to the Neogene had significantly lower CSI values indicating
285 finer substrate types than geologies from the Carboniferous to the Precambrian (Fig. 3).

286 The geology age categories were recoded into one indicator variable coding for solid
287 geology ages younger than the Carboniferous (Fig. 4). Geology age distribution separates
288 Wales, Cornwall and part of the North and the Lake District from the lowland areas of
289 eastern and southern England.

290 Hydrometric areas were recombined into six groups based on their average CSI value
291 and ordered according to increasing substrate size. The distribution of hydrometric area
292 categories follows a pattern separating the uplands from the lowlands of England and
293 Wales (Fig. 5). Hydrometric area categories represent catchment size and its influence on
294 substrate with smaller hydrological units displaying coarser substrate types than larger
295 lowland catchments. Group 4, 5 and 6 represent steeper catchments with higher levels of
296 hill slope activity and reflect the preponderance of upland controls in the delivery of
297 sediments. Catchments from the lower groups tend to have a higher proportion of
298 streams originating and running in low altitude low slope areas compared to higher

299 categories. Coarse sediment tends to originate closer to source and is linked to local
 300 erosion of hard rocks generally located in the upland areas. A lack of upland control within
 301 catchments is, therefore, likely to result in lower delivery of coarse sediments within the
 302 river system and a higher proportion of fine sediment arising from downstream attrition
 303 and fining (Werritty, 1992).

304



305

306 Fig. 3: One way Anova of CSI value against solid geology age categories derived from
 307 the 1979 BGS solid geological map for all RHS sites. Average CSI values per age
 308 category with 95% confidence intervals based on pooled standard deviation ($F=819$,
 309 $p<0.0001$, $n=9934$).

310 Solid geologies were grouped into eight classes based on increasing average CSI value
 311 (Fig. 6). Upon examination, solid geology categories reflect two related aspects:
 312 geological age and hardness. The first class contains recent erodible clay and limestone
 313 formations and displays rivers with predominantly fine sediment material. The following
 314 two solid geology classes comprise slightly older (Jurassic) soft sedimentary formations

315 including chalk, clay, limestone, shales and marls that support streams dominated by fine
316 sediments with some occurrence of gravels and pebbles. Class three is dominated by
317 sedimentary sand, clay and Oolitic geologies. River channels running on those geologies
318 tend to display a higher fraction of gravels and pebbles with a lower predominance of fine
319 sediment. Solid geology class four is constituted of older and harder geologies from the
320 Carboniferous/Triassic period with sandstone and coal that support rivers with a
321 significantly higher occurrence of coarse substrate such as gravel-pebbles and cobbles
322 and little fine sediment. Class five contains geologies from the Cambrian up to the
323 Carboniferous with a predominance of metamorphic and intrusive rocks such as grit stone
324 and granite. Rivers running on these profiles tend to have coarser substrate with cobbles,
325 boulders and bedrock. Class six is constituted mainly of hard igneous and Palaeozoic
326 sedimentary rocks. These are associated with rivers showing a dominance of cobbles
327 with greater occurrence of boulders and bedrock. The last class represents Cambrian grit
328 and limestone rocks that are characterised by very coarse substrate types.

329 Geographically, harder and older geologies are located in the west of England, in Wales,
330 in the North West and near the Scottish border (Fig. 6).

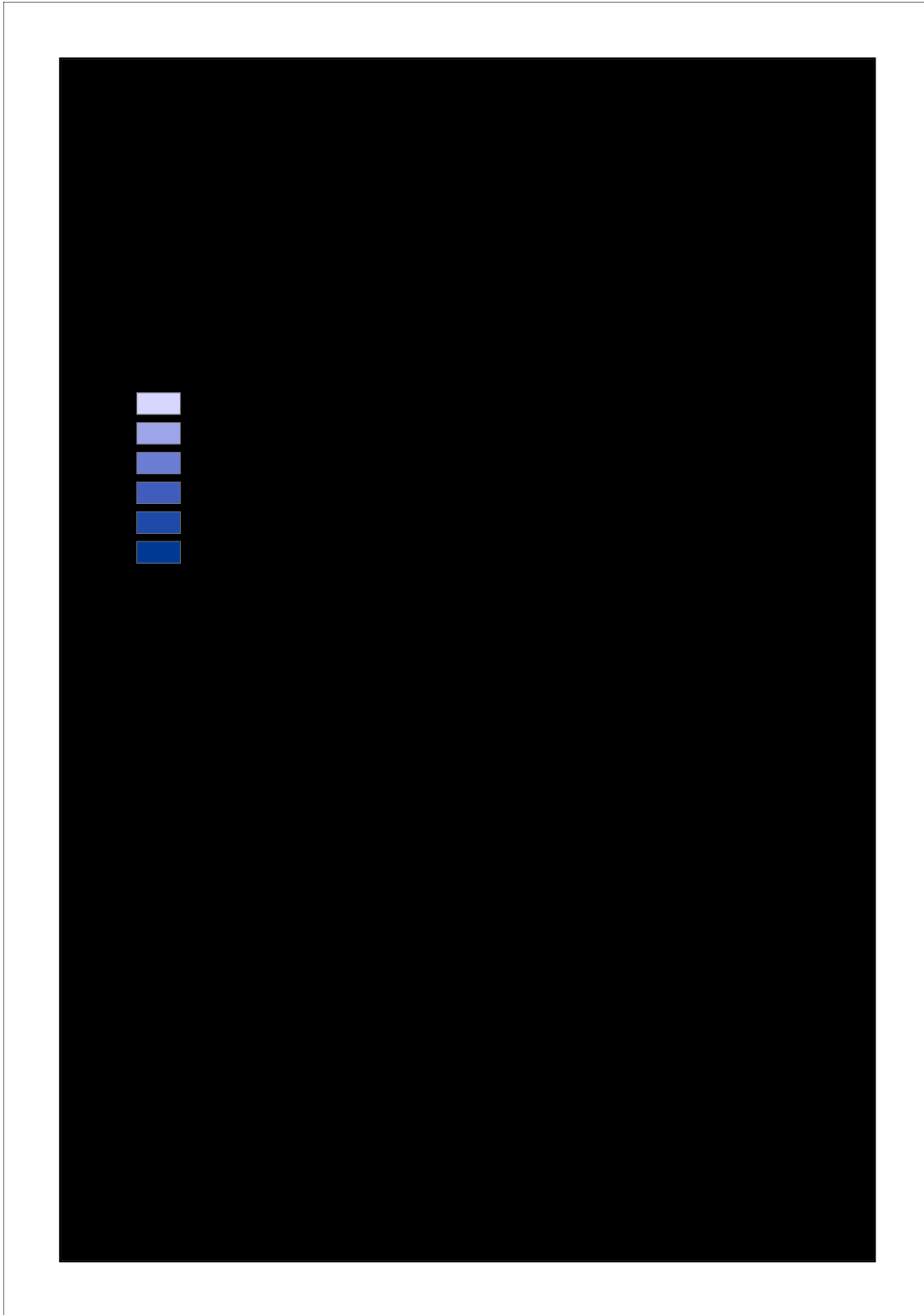
331 From the previous three maps, we can observe spatial correlations between geological
332 age, solid geology classes and hydrometric area categories. Although the three sets of
333 variables are correlated, they each provide subtle differences in explaining substrate
334 distribution.



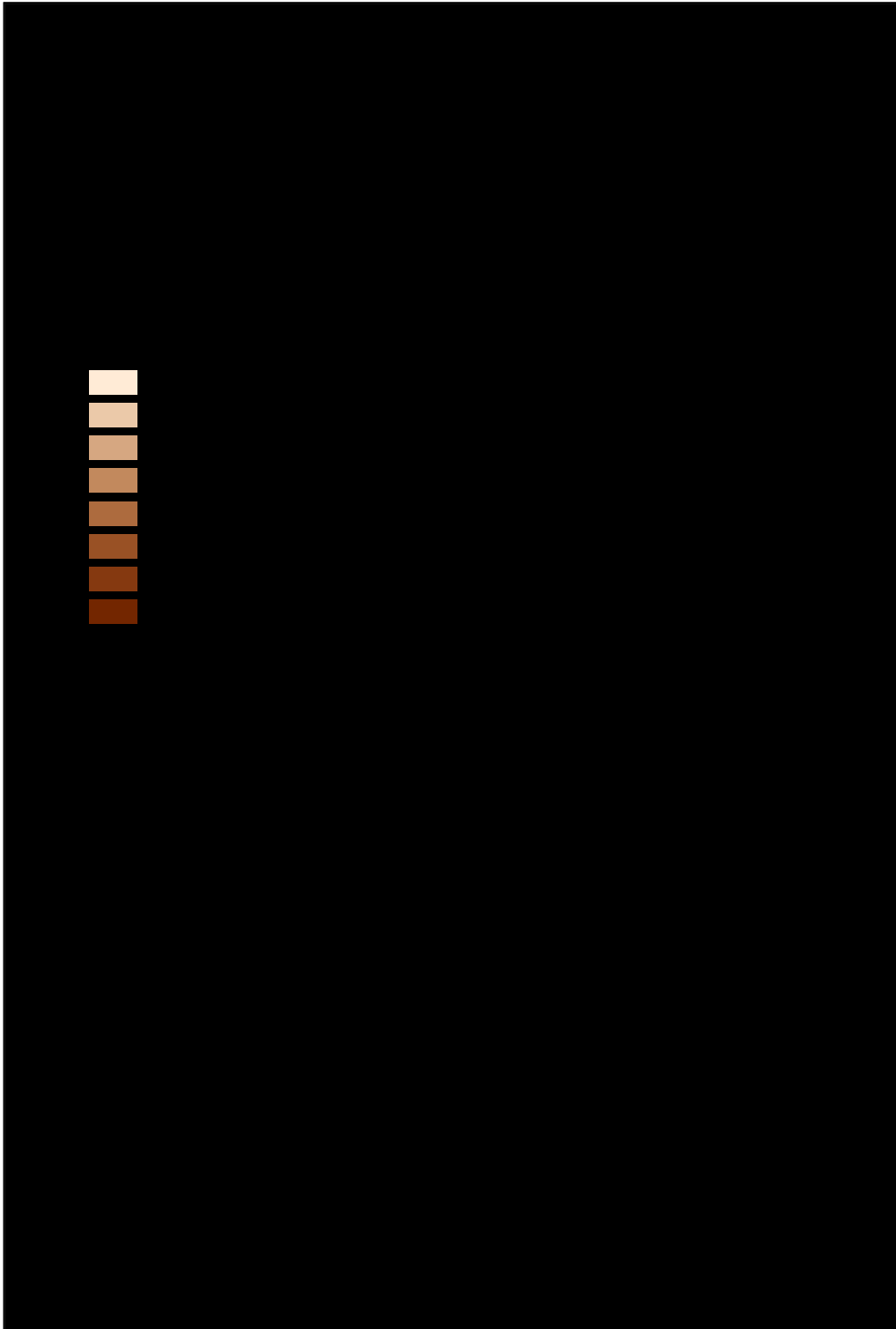
335

336 Fig. 4: Solid geology age distribution in England and Wales from the 1:625,000 BGS solid
337 geology map recombined into two classes. 'Recent geologies' represent geologies from
338 the Neogene to the Triassic and 'older geologies' from the Carboniferous to the
339 Precambrian.

340



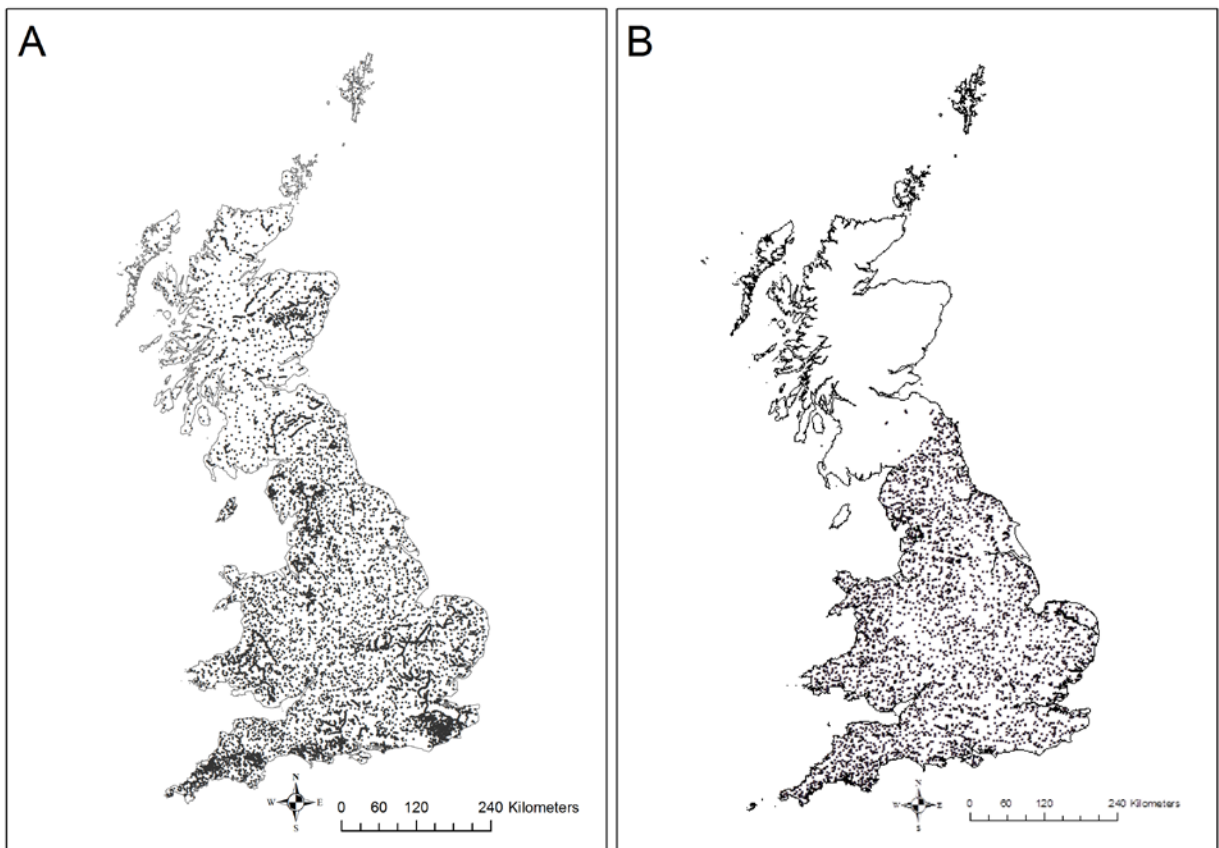
341
342 Fig. 5: Hydrometric area category distribution in England and Wales. Categories from 1
343 to 6 represent hydrological units with increasing average CSI values for surveyed RHS
344 sites. Low CSI values correspond to fine sediment dominated streams, high CSI to
345 coarse sediment dominated streams.



346
347 Fig. 6: Solid geology class distribution in England and Wales. Classes from 0 to 8
348 represent solid geologies with increasing average CSI value for surveyed RHS sites. Low
349 CSI values correspond to fine sediment dominated streams, high CSI to coarse sediment
350 dominated streams.

351 3.3. Regression kriging

352 Following data quality checks, 9473 British RHS sites were retained for the analyses (Fig.
353 7A). The best model, following selection, included the four PCA axes, two solid geology
354 categories, geological age and four hydrometric area groups (Table 2). The model
355 explained 67% of the variability in CSI.



356
357 Fig. 7 : Distribution of A) 9473 RHS sites used for modelling CSI and B) 3884 sites used
358 for testing.

359 The best fit for modelling the OLS and GLS residual variograms was obtained using a
360 combination of spherical and exponential functions (Fig. 8). Spatial correlations were
361 observable up to 5 km and started to plateau after 13 km. The presence of a nugget can
362 be explained by between-surveyor variability as well as time of survey, flow condition etc.
363 The residuals showed a slight tendency for under-prediction, but no great departure from

364 normality (Fig. 9A). Kriged residuals were added to the cross-validated linear model
 365 predictions and the estimated R^2 for the spatially corrected model was 0.74. The kriged
 366 model residuals showed a marked improvement in prediction with a tighter and more
 367 symmetrical distribution around the mean and figures very close to zero (Fig. 9B). A
 368 variogram plot of residuals following kriging showed no sign of remaining spatial
 369 correlation. The kriged model residuals were also investigated for different years of
 370 survey to check the validity of the model over time. The residuals showed no clear pattern
 371 of change in distribution between the years of survey with most variability explained by
 372 differences in sample size.

Table 2: OLS model linking transformed CSI values to GIS map-derived covariates following a best subset selection procedure.

Model Summary

Model	R	R ²	Adjusted R ²	RMSE
1	0.818	0.670	0.669	1.103

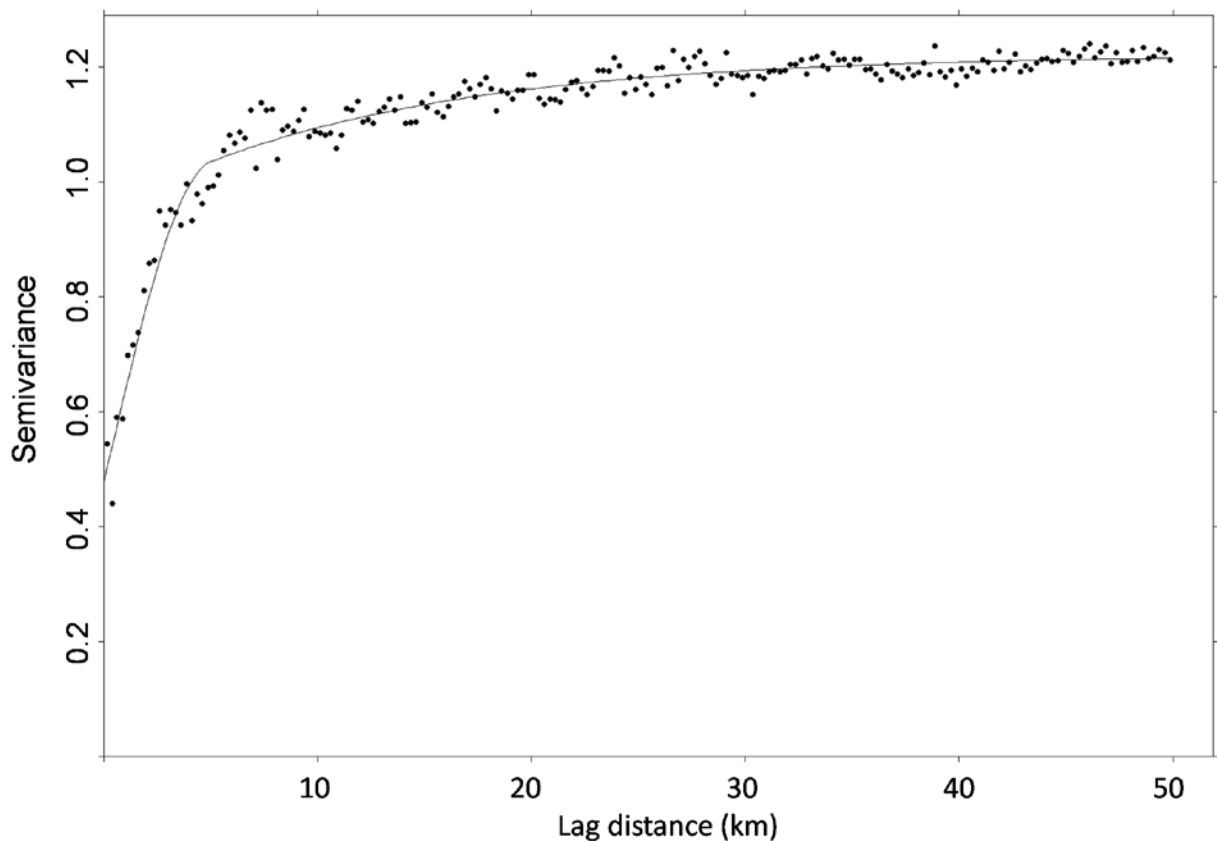
	Sum of Squares	df	Mean Square	F	p
Regression	23345	11	2122.265	1743	< .001
Residual	11522	9466	1.217		
Total	34867	9477			

	Unstandardised	Standard Error	Standardised	t-value	p
intercept	5.978	0.029		209.358	< .001
PCA1	0.566	0.011	0.425	53.725	< .001
PCA2	0.443	0.011	0.262	39.218	< .001
PCA3	0.158	0.021	0.048	7.487	< .001
PCA4	0.062	0.028	0.015	2.240	0.025
Geological age	-0.886	0.048	-0.231	-18.386	< .001
Solid Geology 2	0.301	0.038	0.073	7.831	< .001
Solid Geology 3	0.585	0.072	0.053	8.094	< .001
Hydrometric group 1	-1.525	0.072	-0.158	-21.032	< .001
Hydrometric group 2	-1.073	0.046	-0.272	-23.284	< .001
Hydrometric group 3	-0.558	0.062	-0.068	-9.075	< .001
Hydrometric group 4	-0.540	0.032	-0.124	-16.712	< .001

375

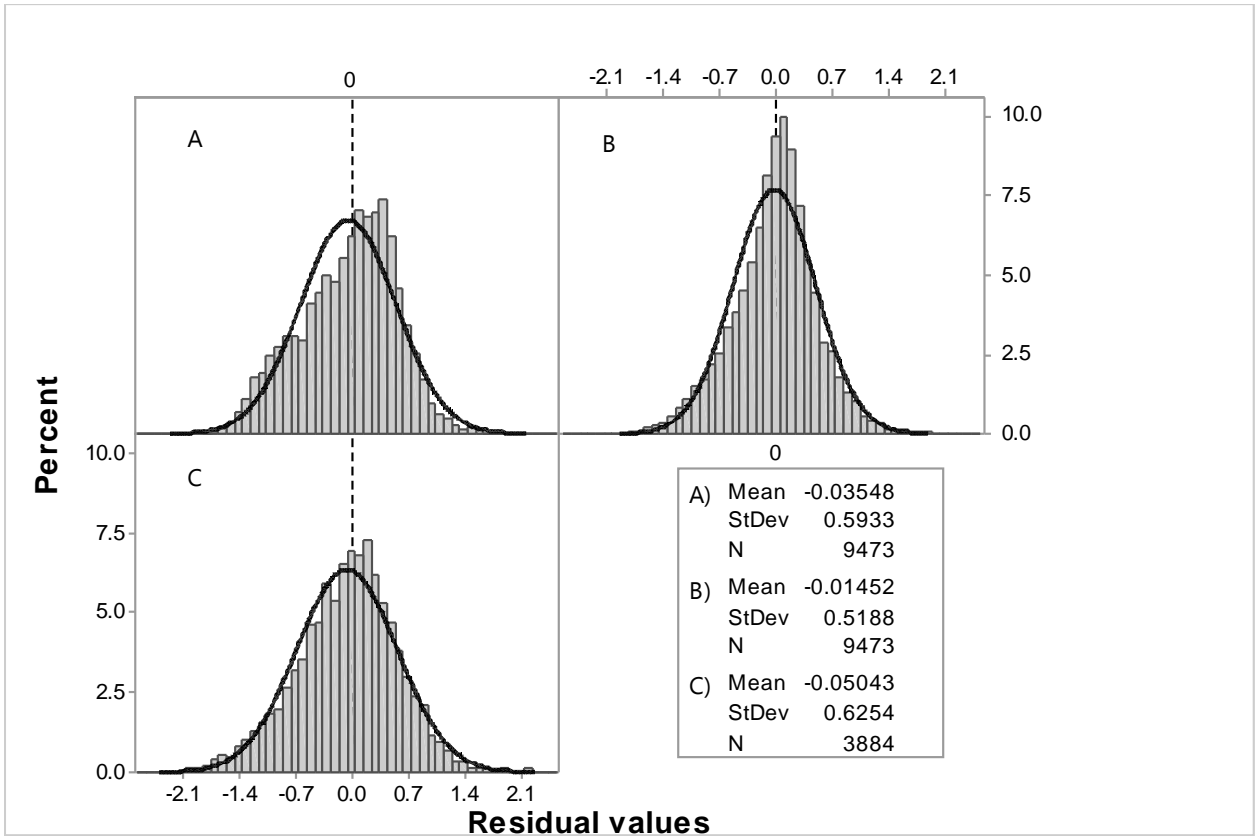
376 The model was then tested on 3884 independent sites surveyed in 2006-11 mainly in
377 England and Wales (Fig. 7B). The model explained 64% of the variability in the new
378 data. The residual distribution was centred around zero and showed no tendency for
379 over- or under-prediction (Fig. 9C). Histograms of residuals per year of survey showed no
380 significant pattern.

381 The model and test data were then joined and the regression kriging model was applied
382 to the entire river network. The resulting map (Fig. 10) shows a clear gradient between
383 the uplands in the West and North of England and Wales dominated by harder, older
384 geologies and coarser substrate types, and the East and South, where sedimentary rocks
385 predominate and channels are dominated by finer sediments.

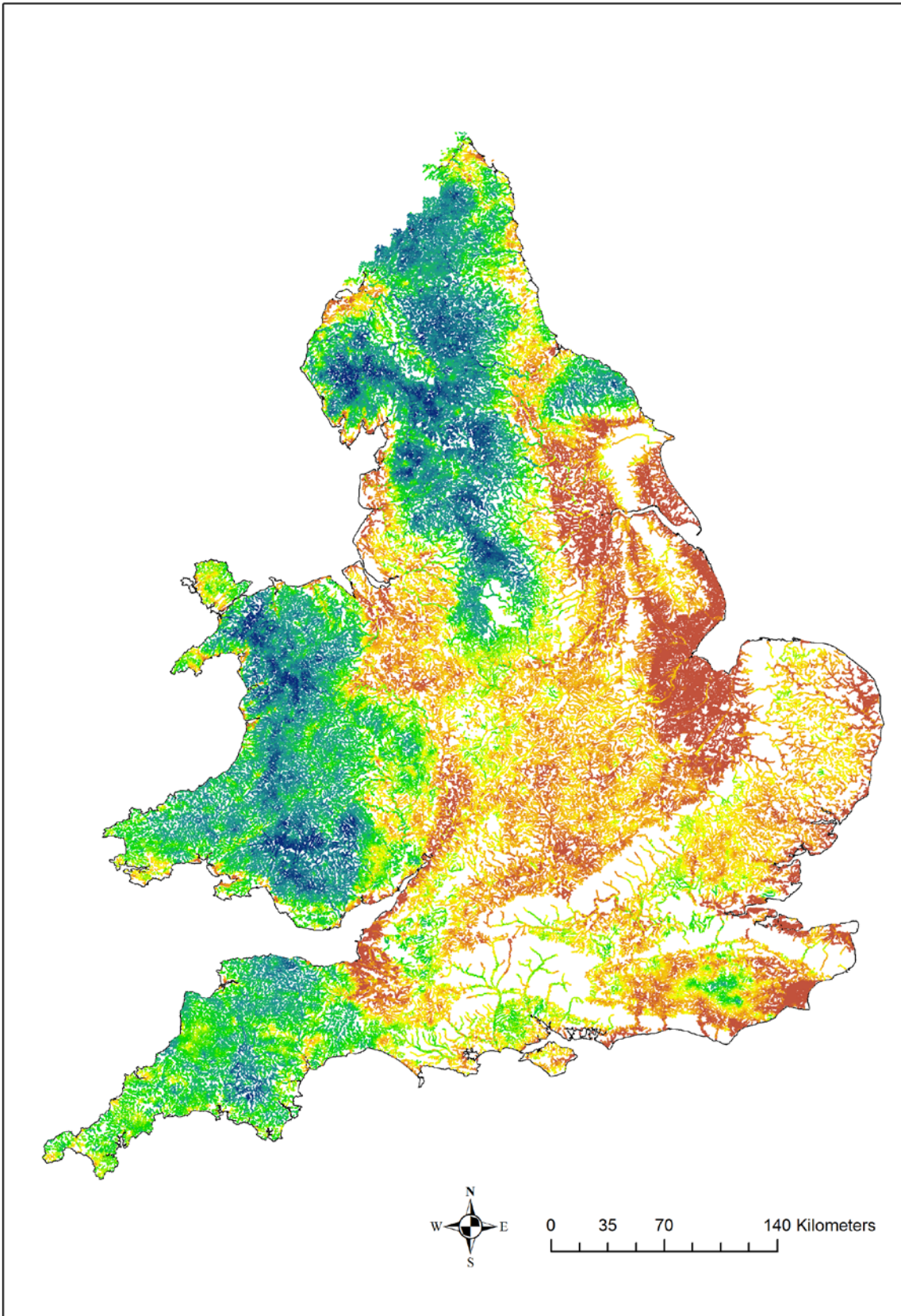


386

387 Fig. 8: OLS model residual variogram fitted with a combination of exponential and
388 spherical functions.



389
 390 Fig. 9: Distribution of model residual values with fitted curves for the cross-validated GLS
 391 linear regression model (A) before and (B) after kriging; and (C) for the test sample.



392
393 Fig. 10: Map of predicted values of CSI using regression kriging at every 500m across the
394 1:50,000 river network on a gradient from bedrock/boulder (blue) to gravel-pebble (green)

395 and silt-sand-clay (brown). White reflects areas of low drainage density where fewer
396 streams are present.

397 **4. Discussion**

398 Using existing data and geostatistical modelling techniques, it was possible to identify and
399 predict channel substrate and apply the model to the entire river network in England and
400 Wales thus, providing environmental practitioners and managers with the first
401 comprehensive national scale map of channel substrate distribution across the network.

402 Traditionally, substrate is characterised using some quantification of sediment size
403 distribution; generally statistics taken from a distribution such as D_{50} or D_{84} (median size
404 or 84th percentile size). Sediment sizing is based on field survey where substrate is either
405 sampled manually or using mechanical techniques (Kondolf et al., 2003). Sampling efforts
406 can be intensive and, therefore, expensive (Bunte et al., 2009) and there are no existing
407 national datasets of substrate size available at present. Davenport et al (2004) attempted
408 to derive estimates of substrate size using RHS data. The Sediment Calibre Index (SCI)
409 was calculated by multiplying substrate occurrence by Wentworth category median size in
410 phi units and averaging over 10 spot-checks. We produced a modified version of the SCI,
411 the SCI_m by adding one category including bedrock and artificial substrate using the
412 following equation and correcting some of the mistakes that were introduced in the
413 original publication for sand and silt phi values (Angela Gurnell, pers. comm.):

$$414 \quad SCI_m = (-12(AR+BE) - 8 BO - 7 CO - 3.5 GP + 1.5 SA + 6 SI + 9 CL) / N_{sc}$$

415 The SCI_m was applied to 10,135 RHS British sites and compared to the CSI. The
416 correlation between the two indices was extremely high (Pearson correlation coefficient =
417 -0.985, $n=10,135$) suggesting that the CSI and the SCI_m both represent average channel
418 substrate size. It is important to note that CSI and SCI_m are unlikely to represent D_{50}

419 unless the substrate size distribution follows a unimodal symmetrical distribution. Further
420 studies involving comparisons with more traditional sediment sizing techniques and RHS
421 would help identify the relevance and significance of the CSI with regards to
422 characterising channel substrate across a 500 m reach but these data currently do not
423 exist.

424 One strong advantage of using multivariate techniques is in the identification of major
425 patterns and dimensions that can be related to biological gradients of species distribution.
426 The assumption is that species distribution and community composition adapt to
427 dominant habitat gradients. In the present case, the main dimension extracted was a well-
428 known substrate fining gradient (Werritty, 1992). Another advantage is that dominant
429 gradients are likely to be influenced by drivers of geomorphological change and are,
430 therefore, more predictable. CSI was correlated and explained by a series of attributes
431 acting at different scales that can be related to known drivers of geomorphological
432 change. At the local (site) scale, Jeffers' PCA1 and PCA3 represent ground slope whilst
433 PCA2 acts as a surrogate for discharge. Slope and discharge are the main drivers of
434 stream power which is strongly related to sediment transport and sorting (Rice and
435 Church, 1998). Solid geology classes and PCA4 provide a wider catchment scale context
436 of geomorphological influence on stream energy and sediment supply. The geology
437 categories reflect the age and hardness of geological types whilst PCA4 represents
438 upstream catchment slope (Jeffers, 1998). Wider scale influences were represented by
439 attributes such as geological age and hydrometric area groups that provided a greater
440 spatial and climatic context for predicting channel substrate.

441 The geostatistical analysis of model residuals revealed the presence of remaining
442 unexplained spatially correlated variance with a relatively short range. This suggests the
443 presence of local random factors influencing substrate distribution, potentially linked to

444 sediment transport, sediment supply from the surrounding landscape (Church, 2002),
445 riparian land use or human-made impacts. They could also represent non-linear spatial
446 responses to geomorphic drivers or local discontinuities in substrate caused by
447 confluences, landslides/bank erosion or the presence of lakes or reservoirs (Rice et al.,
448 2001).

449 The kriging process greatly reduced the spatial correlation in the residuals and increased
450 the model predictive power. The model was also tested for its ability to predict substrate
451 for different surveys. Channel substrate and geomorphological forms tend to be quite
452 stable over decadal timescale and evolve slowly unless significant changes occur such as
453 channel modification or catastrophic events (Knighton, 1998). Therefore, we expected to
454 see no large decrease in predictive power across the years of survey. This was confirmed
455 for both the modelling and the test samples which showed very little deviation in
456 predictive power across survey years.

457 The overall predictive power of the model on the test sample was satisfactory with 64% of
458 the site variability explained by the model. An examination of model residuals for the 100
459 sites with the largest residual values showed no discernible patterns apart from under-
460 prediction of artificial, bedrock (22 % of the sites) and peat substrates (2% of sites). Other
461 sources of error were investigated. For the 2007-8 RHS baseline survey, one third of sites
462 were surveyed on parts of the 1:50,000 river network not covered by the previous
463 sampling strategy which was based on the 1:250,000 river network. A comparison of
464 residual values between sites located on the 1:250,000 and 1:50,000 networks using
465 ANOVA showed a significant difference between sample means ($F=384$; $p<0.001$;
466 $n=2772$). The model tended to predict slightly coarser substrate size on the 1:50,000
467 sample than observed. This could be linked to stream size and management regime.
468 Sites selected on the 1:50,000 sample tended to be narrower, with an average bankfull

469 width of 4.1 m ($n=1692$), compared to 9 m ($n=3134$) for the 1:250,000 sites, with a strong
470 presence of agricultural ditches and artificial channels. It is possible that the model
471 parameters do not fully account for small artificial channels and this shows some of the
472 limitations of using the predictive model on sites collected at different scales. The overall
473 impact of scale on predictive accuracy was, however, small and produced only a slight
474 decrease in overall predictive accuracy. A practical advantage of RK is that it honours
475 field data. Thus, predictions using RK will always fit perfectly the observed values at
476 surveyed sites. This is important from a predictive accuracy viewpoint, but also on an
477 operational viewpoint as it reinforces the credibility of the model in the eyes of users
478 (Naura, 2014).

479 **5. Conclusion**

480 We proposed an alternative approach for mapping habitat elements across *entire* river
481 networks that makes use of existing semi-quantitative survey data, GIS map-based
482 covariate data and RK. A new national scale substrate index has been developed, which
483 is accurate from 500m up to national scales. This application shows the potential power of
484 using spatially explicit techniques for modelling river attributes at the national scale. The
485 analyses presented in this article are part of a broader effort to characterise and map river
486 habitats, identify river reaches for environmental management and develop practical tools
487 for impact assessment, diagnostics and management planning that will be demonstrated
488 in subsequent publications.

489 **Acknowledgements**

490 RHS data can be accessed through the UK government data portal
491 (<http://data.gov.uk/>). The research was partly funded through grants provided by the
492 School of Geography at the University of Southampton, the Environment Agency, the

493 Scottish Environment Protection Agency and the EPSRC. The Authors would like to
494 dedicate this work in memory of Prof Mike Clark, whose dedication to the science and
495 practice of environmental and specifically river management is greatly missed.

496 **References**

- 497 Bivand, R.S., Pebesma, E.J., Gomez-Rubio, V., 2008. Applied spatial data analysis with
498 R. Springer, New York.
- 499 Bunte, K., Abt, S.R., Potyondy, J.P., Swingle, K.W., 2009. Comparison of Three Pebble
500 Count Protocols (EMAP, PIBO, and SFT) in Two Mountain Gravel-Bed Streams. JAWRA
501 Journal of the American Water Resources Association 45, 1209-1227.
- 502 CEN, 2004. A guidance standard for assessing the hydromorphological features of rivers.,
503 in: Comité Européen de Normalisation (Ed.).
- 504 Cherrill, A., Mcclean, C., 1999. Between-observer variation in the application of a
505 standard method of habitat mapping by environmental consultants in the UK. Journal of
506 Applied Ecology 36, 989-1008.
- 507 Chessman, B.C., Fryirs, K.A., Brierley, G.J., 2006. Linking geomorphic character,
508 behaviour and condition to fluvial biodiversity: implications for river management. Aquatic
509 Conservation-Marine and Freshwater Ecosystems 16, 267-288.
- 510 Church, M., 2002. Geomorphic thresholds in riverine landscapes. Freshwater Biology 47,
511 541-557.
- 512 Collen, B., Whitton, F., Dyer, E., Baillie, J., Cumberlidge, N., Darwall, W., Pollock, C.,
513 Richman, N., Soulsby, A., Bohm, M., 2014. Global patterns of freshwater species
514 diversity, threat and endemism. Global Ecology and Biogeography 23, 40-51.
- 515 Collins, A.L., Jones, J.I., Sear, D.A., Naden, P.S., Skirvin, D., Zhang, Y.S., Gooday, R.,
516 Murphy, J., Lee, D., Pattison, I., Foster, I.D.L., Williams, L.J., Arnold, A., Blackburn, J.H.,
517 Duerdoth, C.P., Hawczak, A., Pretty, J.L., Hulin, A., Marius, M.S.T., Smallman, D.,
518 Stringfellow, A., Kemp, P., Hornby, D., Hill, C.T., Naura, M., Brassington, J., 2014.
519 Extending the evidence base on the ecological impacts of fine sediment and developing a

520 framework for targeting mitigation of agricultural sediment losses. Defra report WQ128,
521 Defra ed. Defra.

522 Davenport, A.J., Gurnell, A.M., Armitage, P.D., 2004. Habitat survey and classification of
523 urban rivers. *River Research and Applications* 20, 687-704.

524 Downs, P.W., 1995. River channel classification for channel management purposes.
525 *Changing River Channels*. John Wiley & Sons Ltd., University of Nottingham, UK, 347-
526 367.

527 Environment Agency, 2012. National Risk Assessments for Morphology: Broad-Scale
528 Evaluation. Environment Agency, Bristol, p. 31.

529 European Environment Agency, 2012. European waters - assessment of status and
530 pressures, p. 100.

531 European Union, 2000. Directive 2000/60/EC of the European Parliament and of the
532 Council of 23 October 2000 establishing a framework for Community action in the field of
533 water policy.

534 Fuller, R.M., Smith, G.M., Sanderson, J.M., Hill, R.A., Thomson, A.G., 2002. The UK
535 Land Cover Map 2000: Construction of a parcel-based vector map from satellite images.
536 *Cartographic Journal* 39, 15-25.

537 Gasparini, N.M., Tucker, G.E., Bras, R.L., 1999. Downstream fining through selective
538 particle sorting in an equilibrium drainage network. *Geology* 27, 1079-1082.

539 Greenacre, M.J., 1993. Correspondence analysis in practice. Academic Press Limited,
540 San Diego.

541 Hendry, K., Cragg-Hine, D., 1997. Fisheries Walkover Surveys. A Guidance Manual,
542 Fisheries Technical Manual, Bristol

543 Hornby, D.D., 2010. RivEX, 6.7 ed. <http://www.rivex.co.uk>.

544 Jeffers, J.N.R., 1998. An ordination of river habitats using RHS data. *Aquatic*
545 *Conservation: Marine and Freshwater Ecosystems*. 8, 529-540.

546 Knighton, D., 1998. *Fluvial forms and processes - A new perspective*. Arnold, London.

547 Kondolf, G.M., Lisle, T.E., Wolman, G.M., 2003. Bed sediment measurement, in: Kondolf,
548 G.M., Piégay, H. (Eds.), *Tools in Fluvial Geomorphology*. John Wiley & Sons, pp. 347-
549 395.

550 Loh, J., Green, R., Ricketts, T., Lamoreux, J., Jenkins, M., Kapos, V., Randers, J., 2005.
551 *The Living Planet Index: using species population time series to track trends in*
552 *biodiversity*. *Philosophical Transactions of the Royal Society B-Biological Sciences* 360,
553 289-295.

554 Maddock, I., 1999. The importance of physical habitat assessment for evaluating river
555 health. *Freshwater Biology* 41, 373-391.

556 Minitab, I., 2010. *Minitab 16 Statistical Software*, 16 ed. State College.

557 Montgomery, D.R., Buffington, J.M., 1997. Channel-reach morphology in mountain
558 drainage basins. *Geological Society of America Bulletin* 109, 596-611.

559 Morris, P.H., Williams, D.J., 1999. A worldwide correlation for exponential bed particle
560 size variation in subaerial aqueous flows. *Earth Surface Processes and Landforms* 24,
561 835-847.

562 Mosley, M.P., 1987. The classification and characterization of rivers, in: Richard, K.S.
563 (Ed.), *River Channels: Environment and Process*. Blackwell, Oxford, pp. 295-320.

564 Naiman, R.J., Beechie, T.J., Benda, L.E., Berg, D.R., Bisson, P.A., MacDonald, L.H.,
565 O'Connor, M.D., Olson, P.L., Steel, E.A., 1992. Fundamental elements of ecologically
566 healthy watersheds in the Pacific Northwest coastal ecoregion, *Watershed management*.
567 Springer, pp. 127-188.

568 Naura, M., 2014. Decisions Support Systems. Factors affecting their design and
569 implementation within organisations. Lessons from two case studies. Lambert Academic
570 Publishing, Berlin.

571 Newson, M.D., Clark, M.J., Sear, D.A., Brookes, A., 1998. The geomorphological basis
572 for classifying rivers. *Aquatic Conservation-Marine and Freshwater Ecosystems* 8, 415-
573 430.

574 Raven, P.J., Fox, P., Everard, M., Holmes, N.T.H., Dawson, F.H., 1997. River habitat
575 survey: A new system for classifying rivers according to their habitat quality. *Freshwater*
576 *Quality: Defining the Indefinable?*, 215-234.

577 Revenga, C., Campbell, I., Abell, R., de Villiers, P., Bryer, M., 2005. Prospects for
578 monitoring freshwater ecosystems towards the 2010 targets. *Philosophical Transactions*
579 *of the Royal Society B-Biological Sciences* 360, 397-413.

580 Rice, S., Church, M., 1998. Grain size along two gravel-bed rivers: Statistical variation,
581 spatial pattern and sedimentary links. *Earth Surface Processes and Landforms* 23, 345-
582 363.

583 Rice, S.P., Greenwood, M.T., Joyce, C.B., 2001. Tributaries, sediment sources, and the
584 longitudinal organisation of macroinvertebrate fauna along river systems. *Canadian*
585 *Journal of Fisheries and Aquatic Sciences* 58, 824-840.

586 Rosgen, D.L., 1994. A Classification of Natural Rivers. *Catena* 22, 169-199.

587 Sala, O., Chapin, F., Armesto, J., Berlow, E., Bloomfield, J., Dirzo, R., Huber-Sanwald, E.,
588 Huenneke, L., Jackson, R., Kinzig, A., Leemans, R., Lodge, D., Mooney, H., Oesterheld,
589 M., Poff, N., Sykes, M., Walker, B., Walker, M., Wall, D., 2000. Biodiversity - Global
590 biodiversity scenarios for the year 2100. *Science* 287, 1770-1774.

591 Sear, D., Newson, M., Hill, C., Old, J., Branson, J., 2009. A method for applying fluvial
592 geomorphology in support of catchment-scale river restoration planning. *Aquatic
593 Conservation-Marine and Freshwater Ecosystems* 19, 506-519.

594 Strayer, D., Dudgeon, D., 2010. Freshwater biodiversity conservation: recent progress
595 and future challenges. *Journal of the North American Benthological Society* 29, 344-358.

596 Thorne, C.R., 1997. Channel types and morphological classification. *Applied fluvial
597 geomorphology for river engineering and management*, 175-222.

598 Tockner, K., Stanford, J.A., 2002. Riverine floodplains: present state and future trend.
599 *Environmental Conservation* 29, 308-330.

600 Townsend, C.R., Hildrew, A.G., 1994. Species traits in relation to a habitat templet for
601 river systems. *Freshwater Biology* 31, 265-275.

602 Vaughan, I., Merrix-Jones, F., Constantine, J., 2013. Successful predictions of river
603 characteristics across England and Wales based on ordination. *Geomorphology* 194,
604 121-131.

605 Vaughan, I.P., Diamond, M., Gurnell, A.M., Hall, K.A., Jenkins, A., Milner, N.J., Naylor,
606 L.A., Sear, D.A., Woodward, G., Ormerod, S.J., 2009. Integrating ecology with
607 hydromorphology: a priority for river science and management. *Aquatic Conservation-
608 Marine and Freshwater Ecosystems* 19, 113-125.

609 Vaughan, I.P., Ormerod, S.J., 2003. Improving the quality of distribution models for
610 conservation by addressing shortcomings in the field collection of training data.
611 *Conservation Biology* 17, 1601-1611.

612 Vorosmarty, C.J., McIntyre, P.B., Gessner, M.O., Dudgeon, D., Prusevich, A., Green, P.,
613 Glidden, S., Bunn, S.E., Sullivan, C.A., Liermann, C.R., Davies, P.M., 2010. Global
614 threats to human water security and river biodiversity. 467, 555-561.

615 Webster, R., Oliver, M.A., 2007. Geostatistics for environmental scientists, 2nd edition ed.
616 John Wiley & Sons Ltd, Chichester.

617 Wentworth, C.K., 1922. A scale of grade and class terms for clastic sediments. The
618 Journal of Geology, 377-392.

619 Werritty, A., 1992. Downstream fining in a gravel bed river in Southern Poland:
620 Lithological controls and the role of abrasion, in: P., B., R.D., H., C.R., T., P., T. (Eds.),
621 Dynamics of Gravel Bed Rivers. Wiley: Chichester, pp. 333–346.

622 World Conservation Monitoring Centre, 1998. Freshwater Biodiversity: a preliminary
623 global assessment, WCMC Biodiversity Series. World Conservation Monitoring Centre,
624 Cambridge, UK, p. 132.

625 WWF, 2014. Living Planet Report 2014. Species and spaces, people and places. WWF,
626 p. 176.

627