

1 Can a paleo-drought record be used to  
2 reconstruct streamflow? A case-study for  
3 the Missouri River Basin

4 Michelle Ho<sup>1</sup>, Upmanu Lall<sup>1,2</sup>, Edward R. Cook<sup>3</sup>

5 1. Columbia Water Center, Columbia University, New York, NY, USA

6 2. Department of Earth and Environmental Engineering, Columbia University, New York, New York, USA

7 3. Lamont-Doherty Earth Observatory of Columbia University, Palisades, New York 10964, USA

8 Corresponding author: Michelle Ho. Email: mh3538@columbia.edu, Phone: 212-854-7081

9

10 Keywords: paleoclimate, regularized canonical correlation analysis, streamflow reconstruction

11 Key points

- 12 • Streamflow is reconstructed from an existing paleoclimate drought index
- 13 • Methodological innovation using rCCA addresses very high dimensional dataset
- 14 • Reconstructed streamflow provides insights into extreme, persistent high and low flow

## 15 Abstract

16 Recent advances in paleoclimatology have revealed dramatic long-term hydro-climatic variations that  
17 provide a context for limited historical records. A notable dataset derived from a relatively dense  
18 network of paleoclimate proxy records in North America is the Living Blended Drought Atlas (LBDA): a  
19 gridded tree-ring based reconstruction of summer Palmer Drought Severity Index. This index has been  
20 used to assess North American drought frequency, persistence and spatial extent over the past two  
21 millennia. Here, we explore whether the LBDA can be used to reconstruct annual streamflow. Relative to  
22 streamflow reconstructions that use tree rings within the river basin of interest, the use of a gridded  
23 proxy poses a novel challenge. The gridded series have high spatial correlation, since they rely on tree  
24 rings over a common radius of influence. A novel algorithm for reconstructing streamflow using  
25 regularized canonical regression and inputs of local and global covariates is developed and applied over  
26 the Missouri River Basin, as a test case. Effectiveness in reconstruction is demonstrated with  
27 reconstructions showing periods where streamflow deficits may have been more severe than during  
28 recent droughts (e.g. the Civil War, Dust Bowl and 1950s droughts). The maximum persistence of  
29 droughts and floods over the past 500 years far exceed those observed in the instrumental record and  
30 periods of multi-decadal variability in the 1500s and 1600s are detected. Challenges for an extension to  
31 a national streamflow reconstruction or applications using other gridded paleoclimate datasets such as  
32 adequate spatial coverage of streamflow and applicability of annual reconstructions are discussed.

## 33 1 Introduction

34 The concern with anthropogenic climate change and its hydrologic impacts has focused interest on how  
35 long term climate variability may impact streamflow (e.g. Nohara et al., 2006; Seager et al., 2007). The  
36 consequence of using short records to “over-allocate” the flows of major rivers are often cited as an  
37 example of the need for long records that can better inform the possible range of long-term variations  
38 of streamflow (Tootle and Piechota, 2006; McGowan et al., 2009; Woodhouse et al., 2010). Continuous  
39 records of streamflow in the US span several decades at best. Advances in paleoclimatology in the past  
40 few decades have provided opportunities across the world to extend the range of hydroclimatic  
41 variability (e.g. Quinn, 1992; Jones and Mann, 2004; Tierney et al., 2010; Gallant and Gergis, 2011; Vance  
42 et al., 2012; Cook et al., 2013b; Devineni et al., 2013; Ho et al., 2015). While considerable uncertainty  
43 clouds the projections of hydroclimatic states towards the end of the 21<sup>st</sup> century, in the near to  
44 medium term paleoclimate information may be crucial to inform the interannual to decadal variability of  
45 regional water availability as indicated by streamflow for reservoir operation, and agricultural and other  
46 water use decisions.

47 Paleoclimate reconstructions have been developed using proxies that typically span the past 1000-2000  
48 years (also known as the Common Era). The North American region has a relatively dense network of  
49 high-resolution paleoclimate proxy records, primarily comprised of tree-ring chronologies. Tree-ring-  
50 proxy records have been used to assess various components of environmental variations (Fritts, 1976)  
51 including drought severity (e.g. Routson et al., 2011; Cook et al., 2015a), pluvials (e.g. Woodhouse et al.,  
52 2005; Pederson et al., 2012), streamflow variability (e.g. Woodhouse et al., 2006; Prairie et al., 2008;  
53 Allen et al., 2013; Devineni et al., 2013), and precipitation frequency (Woodhouse and Meko, 1997) in  
54 addition to enabling comparisons of past climate with projected climate scenarios (e.g. Ault et al., 2014;  
55 Cook et al., 2015a; Smerdon et al., 2015).

56 Studies focused on the reconstruction of tree-ring-based paleo-hydrology (e.g. annual and season  
57 streamflow and floods) typically utilize proxies derived from tree-ring networks within or near the  
58 catchment region as predictors (e.g. Woodhouse et al., 2006; St. George, 2010; Pederson et al., 2012;  
59 Devineni et al., 2013). However, these networks are spatially irregular with record lengths varying across  
60 chronologies. An alternative to using spatially and temporally irregular tree-ring chronologies as model  
61 predictors is to use an existing derivative of these records, namely the Living Blended Drought Atlas  
62 (LBDA) (Cook et al., 2010a). The LBDA is a paleoclimate reconstruction of the summer (June-August)  
63 Palmer Drought Severity Index (PDSI) that is gridded across North America on a  $0.5^\circ \times 0.5^\circ$   
64 latitude/longitude grid with reconstructions dating back as far as 2000 years. These records are  
65 temporally complete over the Conterminous United States (CONUS) from 1473 onward. The LBDA, or its  
66 predecessor the North American Drought Atlas (Cook and Krusic, 2004), have been used to assess the  
67 frequency and spatial distribution of droughts over the past millennia (e.g. Herweijer et al., 2007; Cook  
68 et al., 2013a).

## 69 2 A proposal for using PDSI to reconstruct streamflow

70 The intent of the modeling case study presented here is to develop a suitable framework with which  
71 streamflow within the CONUS may be reconstructed using a tree-ring-based reconstruction of the PDSI.  
72 In developing the modelling framework we consider: 1) the constraints posed by the LBDA and  
73 implications for reconstructing streamflow; 2) possible temporal resolutions (monthly, seasonal or  
74 annual) for direct streamflow reconstruction using the LBDA data; 3) how to best use local and far-field  
75 LBDA information for local streamflow reconstruction; and 4) provides insights from the 500 year  
76 reconstruction of the multi-site Missouri River Basin flows as to the decadal and longer variability of  
77 streamflow in the region. Given the existence of the spatially and temporally complete LBDA record over  
78 the last 500 years covering the CONUS, we explore whether the LBDA, a reconstruction of PDSI using

79 tree-ring chronologies, could be a reasonable predictor of streamflow variability. In this case the variable  
80 to be reconstructed is annual streamflow with the aim of eventually reconstructing paleoclimate records  
81 of streamflow across the CONUS, an undertaking that has not previously been attempted using the  
82 LBDA or tree-rings.

83 The motivation for using the LBDA stems from our understanding that the growth of moisture-limited  
84 trees, from which the LBDA is derived, are in part governed by climatic forcings that drive soil moisture  
85 availability. That is, given a vector of unspecified climate variables,  $\mathbf{C}_t$  (where  $\mathbf{C}$  may be comprised of, but  
86 not limited to, climate variables such as temperature, rainfall, wind, soil moisture and radiation) we can  
87 define PDSI as  $PDSI_t = f_1(\mathbf{C}_t)$ . Streamflow is also a derivative of a number of climate variables and can  
88 similarly be defined as  $Q_t = f_2(\mathbf{C}_t)$ . We seek to determine if it is possible to derive and fit a function  $f_3$  that  
89 relates streamflow to PDSI where  $Q_t = f_3(PDSI_t)$ . Where suitable instrumental records of climate are  
90 available, we may derive  $f_1$  and this has been approximated in part using methods such as the  
91 Thornthwaite potential evapotranspiration (PET) equation (Thornthwaite, 1948) and the Penman-  
92 Monteith PET equation (Monteith, 1965) to varying degrees of success (Lockwood, 1999; Sheffield et al.,  
93 2012; Dai, 2013). Consideration of the joint probability distributions  $f(PDSI_t, \mathbf{C}_t)$  and  $f(Q_t, PDSI_t)$  enables  
94 the vector of climate variables to be reconstructed by capitalizing on the climate-PDSI relationship and  
95 the tree-ring chronology-PDSI relationship given  $f(\mathbf{C}_t | PDSI_t) f(PDSI_t | \textit{tree-ring chronology}_t)$ . Paleoclimate  
96 streamflow could similarly be derived by implementing  $f_2$ , namely  $f(Q_t | \mathbf{C}_t) f(\mathbf{C}_t | PDSI_t) f(PDSI_t | \textit{tree-ring}$   
97 ***chronology}\_t). However, the challenge in this approach is the selection of an appropriate vector of  
98 climate variables, many of which may be sparsely observed or unobserved in the instrumental record.  
99 Therefore, we consider modelling paleoclimate streamflow through  $f(Q_t | PDSI_t) f(PDSI_t | \textit{tree-ring}$   
100 ***chronology}\_t). Since the reconstructed PDSI<sub>t</sub> in the LBDA is really  $E[PDSI_t | \textit{tree-ring chronology}_t]$ , we start  
101 by considering  $f(Q_t | PDSI_t)$ , where  $Q_t$  is the streamflow at one or more locations in a river basin, and  
102  $PDSI_t$  represents a vector of LBDA values at the gridded locations of the LBDA that can be a potential******

103 predictor of the streamflows. The relationship can be developed using contemporaneous values of  $PDSI_t$   
104 and historical  $Q_t$ . Subsequently, we can apply this relationship to the paleo estimates of  $PDSI_t$   
105 recognizing that we could use the expected values of  $PDSI_t$  reported in the LBDA, or simulations from  
106 the uncertainty distributions of  $PDSI_t$  reported in the LBDA, as conditioning variables to derive  
107 sequences of  $Q_t$ .

108 A key motivation for using the LBDA is that it is spatially complete across CONUS over the past 500 years  
109 and a successful streamflow reconstruction would have significant value in assessing national water  
110 planning and use strategies and in investigating the different temporal and spatial structures in  
111 streamflow, which differ from the LBDA (see Figures S1 and S2 in supporting information). Relative to  
112 reconstructions that use tree rings within the river basin of interest, the use of a gridded proxy series to  
113 reconstruct watershed processes such as streamflow poses a novel challenge. The gridded LBDA series  
114 have high spatial correlation, since they rely on tree rings over a common radius of influence around  
115 each grid point, nominally 450 km in this case or roughly the correlation decay  $e$ -folding distance  
116 between grid points of the instrumental PDSI. A novel algorithm for reconstruction that uses local and  
117 global covariates for streamflow reconstruction using regularized canonical regression is developed and  
118 applied over the Missouri River Basin, as a test case. This provides a proof of concept at a sub-  
119 continental scale as to whether the approach is feasible. The Missouri River Basin was chosen as it  
120 contains the only major river headwaters in the western US where extensive reconstructions of  
121 paleoclimate hydrology have not been undertaken (Driscoll, 2013) and also parallels current efforts to  
122 further develop tree-ring chronologies and streamflow reconstruction models using tree rings in the  
123 region (Pederson, 2013). A description of the case study region and data are presented in the following  
124 section while Section 4 presents initial diagnostics that inform the modeling approach developed in  
125 Section 5. Section 5 details how the above proposal of using PDSI to reconstruct streamflow is  
126 implemented in the Missouri River Basin. Section 6 provides model verification results and summaries of

127 the key modes of variability in the 500 year mean streamflow reconstructions of Missouri River Basin  
128 streamflow. The final section reviews outstanding questions as to modeling uncertainties, and the  
129 challenges for extending the model presented here to a national scale.

## 130 3 Case Study Region and Data

### 131 3.1 The Missouri River Basin

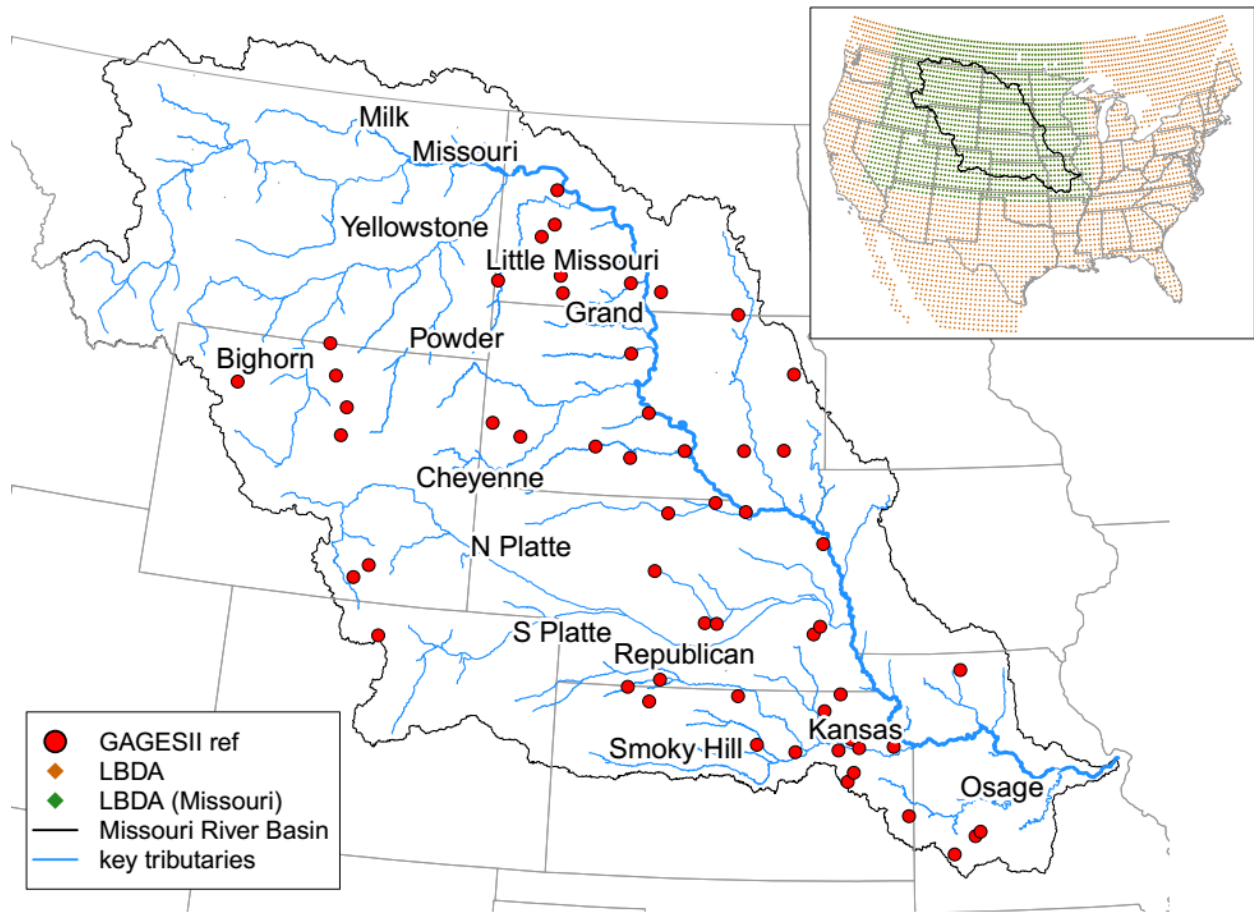
132 The Missouri River is the longest river and the second largest river basin in the US, draining an area over  
133 1.3 million km<sup>2</sup> that spans the southern portions of two Canadian provinces and ten states in the US  
134 (Missouri River basin boundary and key tributaries shown in Figure 1). The headwaters, which are largely  
135 snow-melt fed, are located in the Northern Rocky Mountains. These waters then flow through a largely  
136 semi-arid region to its confluence with the Mississippi River near St. Louis, Missouri (Galat et al., 2005).  
137 Land use in the Missouri River Basin is dominated by agricultural activities including cropping and  
138 grazing which cover 95% of the region (US Army Corps of Engineers, 2006), while the remaining land is  
139 used for recreation, transport, urban and industrial use including mining and energy sector activities  
140 (Galat et al., 2005). An improved perspective of streamflow variability would therefore be of benefit to  
141 the region in terms of managing and balancing the demands for water amongst the various sectors and  
142 users.

143 Precipitation and moisture availability in the Missouri River Basin are characterized by high precipitation  
144 in the western mountainous region, which averages over 1000 mm/year to the drier region in the rain  
145 shadow east of the Rocky Mountains where average annual precipitation is less than 400 mm/year.  
146 Precipitation increases towards the far eastern regions of the Missouri River Basin (Kunkel et al., 2013).  
147 Winter precipitation in the northern Missouri River Basin and in the mountainous regions to the west is  
148 related to the El Niño Southern Oscillation (ENSO) signal from the preceding summer and autumn

149 (Redmond and Koch, 1991). El Niño teleconnections typically manifest as upper level anticyclonic high  
150 pressure cells over the northwestern US and result in the northward displacement or splitting of the jet  
151 stream and anomalously dry conditions in the Missouri River Basin (Trenberth et al., 1988; Dettinger et  
152 al., 1998; Smith et al., 1998). Conversely, La Niña events typically result in wetter conditions in this  
153 region. ENSO impacts are modulated by the decadal scale variability in the northern Pacific Ocean with  
154 warm decadal phases resulting in a deep Aleutian low and corresponding ridging over the western US  
155 thereby enhancing El Niño conditions (Gershunov and Barnett, 1998; Brown and Comrie, 2004). The  
156 enhancement of La Niña impacts by a cool phase in the northern Pacific decadal signature is particularly  
157 noticeable in the northern Missouri River Basin (Wise, 2010). Both Pacific and Atlantic Ocean influences  
158 are seen in the Northern Great Plains with the Great Plains low level jet, originating from the Gulf of  
159 Mexico, enhancing summer precipitation in this region (Higgins et al., 1997).

## 160 3.2 Streamflow data

161 Monthly streamflow data for the Missouri River Basin were obtained from the United States Geological  
162 Survey (USGS) Surface-Water Daily Data for the Nation (<http://nwis.waterdata.usgs.gov/nwis/sw>). The  
163 streamflow data are from stations included in the USGS's GAGES-II network (U.S. Geological Survey,  
164 2011) within the Missouri River Basin boundary and are reference gauges identified by the USGS as the  
165 least-disturbed watersheds with minimal regulation. The selected gauges meet a criterion of data  
166 spanning 40 years with less than 5% missing data (Figure 1, Table S1) and results in 55 streamflow  
167 gauges, 46 of which are also in the USGS Hydro Climatic Data Network.



168

169 *Figure 1. The Missouri River Basin, key tributaries, locations of gauged streamflow used in the analysis, and (inset) the LBDA*  
 170 *grids used in the analysis. LBDA grids in the Missouri River Basin region shown in green.*

171 Three of the selected stations had missing monthly values. These values were imputed using multiple  
 172 imputation by chained equations (MICE) and a method of predictive mean matching (Buuren and  
 173 Groothuis-Oudshoorn, 2011). Monthly streamflow imputation was conducted using streamflow from  
 174 the three closest stations and a cosine function to represent a seasonal signal. The number of  
 175 repetitions in MICE was determined using a rule of thumb method proposed by White et al. (2011) and  
 176 the multiple imputed monthly values were averaged across the repetitions.

177 Monthly data was aggregated into annual streamflow data using a calendar year instead of a water year  
 178 because the average driest month of streamflow at most stations occurred in either December or  
 179 January. Start and end years with incomplete data were excluded. The resulting annual data spans from

180 1929 to 2014 with annual record lengths varying between 39 and 85 years after aggregation. Streamflow  
181 was logarithmically transformed since that leads to a nearly Gaussian distribution for annual streamflow.  
182 Annual streamflow records of zero were replaced with half the minimum annual streamflow prior to the  
183 application of the log transform.

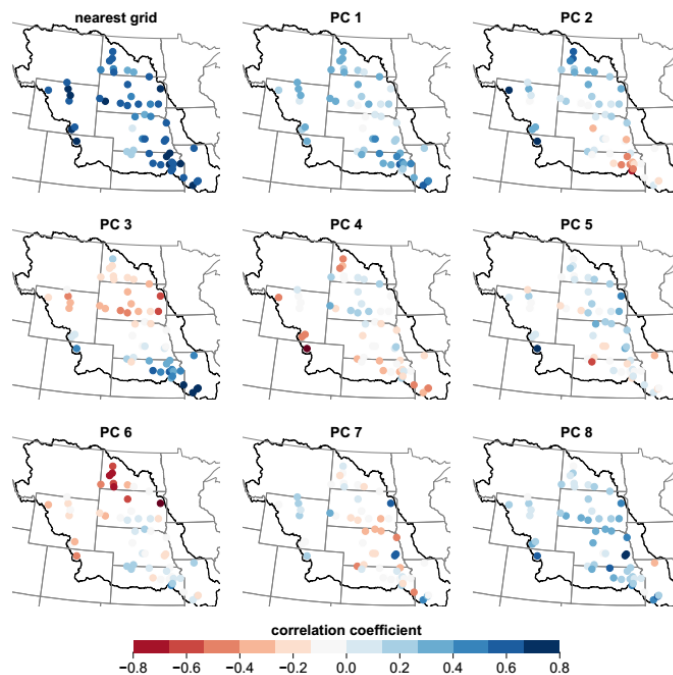
### 184 3.3 A Paleoclimate Record of North American Drought

185 The LBDA is an updated version of the seminal North American Drought Atlas (NADA) (Cook et al., 1999;  
186 Cook and Krusic, 2004; Cook et al., 2010a), which is a paleoclimate reconstruction of the summer (June  
187 to August – JJA) PDSI based on a network of tree-ring chronologies. The LBDA has a spatial resolution of  
188  $0.5^\circ \times 0.5^\circ$  latitude/longitude and incorporates information from 1845 tree-ring chronologies, an  
189 improvement over previous NADA versions that were informed by fewer tree-ring chronologies and  
190 were calculated over coarser grids (Cook et al., 1999; Cook et al., 2004). The LBDA is spatially complete  
191 over the CONUS region from 1473-2005 and includes instrumental data from 1979 onwards. A region of  
192 the LBDA ranging from  $23^\circ\text{N}$  to  $52^\circ\text{N}$  and  $125^\circ\text{W}$  to  $66^\circ\text{W}$  (see green dots in Figure 1) was extracted for  
193 the analysis. This broad region extending beyond the CONUS region was selected to capture patterns of  
194 LBDA variability relevant to the CONUS. A comparative analysis between an instrumental-based gridded  
195 PDSI dataset and the LBDA showed that they are highly correlated, with correlation values significant at  
196 the 99% level and similar variance (results not shown here). The instrumental LBDA data post-1978 were  
197 therefore also included in the modeling analysis performed here.

## 198 4 Initial Diagnostic Analyses

199 Initial diagnostic analyses were performed using both parametric and non-parametric correlations  
200 (Pearson and Kendall correlations respectively) between the LBDA and the log-normal streamflow  
201 series. The two correlations measures yielded similar results suggesting a near-Gaussian linear  
202 dependency (Pizarro and Lall, 2002). Different levels of temporal aggregation were tested including

203 monthly, rolling seasonal, bi-seasonal and annual streamflow. Different representations of the LBDA  
 204 were also tested including using the LBDA grid located nearest to each streamflow gauge, using LBDA  
 205 grids surrounding each streamflow gauge within a given diameter, and principal components (PCs)  
 206 (Jolliffe, 2002) and archetype analysis (Cutler and Breiman, 1994) of US-wide and Missouri River Basin  
 207 region LBDA (see orange and green diamonds respectively in Figure 1). An annual temporal aggregation  
 208 of streamflow was found to produce the strongest signal (Figure 2 showing correlation results between  
 209 streamflow gauges in the Missouri River Basin and the nearest LBDA grid and PCs of US-wide LBDA).  
 210 Diagnostic tests using the nearest LBDA grid were superior and reflects the ability of the point-by-point  
 211 regression method used for the LBDA to preserve local climate details in the PDSI reconstructions that  
 212 are also related to streamflow.



213  
 214 *Figure 2. Pearson correlation between annual streamflow and a) closest LBDA grid b)-i) PC1 – PC8 of US-wide LBDA*  
 215 High correlations were also noted for some streamflow stations with LBDA in the surrounding region  
 216 (see supporting information) particularly for streamflow records with weaker correlations with the  
 217 nearest grid (e.g. gauges in south-central North Dakota and on the border of Nebraska and Kansas). We

218 therefore considered information from LBDA grids within a 450 km radius consistent with the tree  
219 chronology search radius used to form the LBDA. Furthermore, given that large-scale climate and  
220 weather patterns influence local hydroclimatic conditions (Woodhouse et al., 2002), one needs to also  
221 consider the relationships with these larger modes of variability.

222 Correlations between streamflow and the first eight PCs of LBDA grids across the US (LBDA locations  
223 shown in Figure 1, correlations in Figure 2 and PC loading patterns shown in Figure 3) show that the  
224 Missouri River Basin streamflow is correlated with PC1, a PC representing overall US LBDA variability.  
225 However, these correlations are weak in comparison with correlations with the nearest LBDA grid. The  
226 north/south loading pattern of PC 2 is reflected in the change in correlation sign for stations north and  
227 south of the Nebraska and Kansas border. Similarly, the east-west difference in the loadings of PC3 leads  
228 to strong positive correlations in the downstream reaches of the Basin in Kansas and Missouri and  
229 negative correlations in Nebraska and Wyoming. The correlation results between streamflow and PCs of  
230 CONUS LBDA suggest that the large-scale modes of variability could provide additional information for  
231 modelling streamflow. Further details of the eight PCs of the CONUS LBDA and selection methods are  
232 provided in the supporting information.

## 233 5 Modeling approach and performance metrics

234 Here, we present a suite of plausible methods that could be implemented to model streamflow using  
235 the LBDA as the covariate(s) (Section 5.1). We justify the selection of a model that incorporates the use  
236 of regularized canonical correlation (rCCA), which is further described in Section 5.2. A description of the  
237 model performance metrics is provided in Section 5.3.

## 238 5.1 Model design and preliminary model assessment

239 We seek to use log-linear models to quantify the relationship between streamflow at individual gauges  
240 and a suite of site-specific LBDA records to facilitate a paleoclimate reconstruction of streamflow in the  
241 Missouri River Basin. Keeping in mind the application to streamflow record extension using the  
242 reconstructed PDSI values available in the LBDA we consider the following steps:

243 1. *Instrumental Period* (32-77 years at 55 locations): Use LBDA reconstructed PDSI to estimate the  
244 relationship between log transformed annual streamflow and a selection of LBDA records specific to the  
245 target streamflow site, namely  $f(\ln(Q_t) | PDSI_t)$ , where  $PDSI_t$  is really the expected value of PDSI informed  
246 by the tree-ring chronologies,  $E[PDSI_t | tree-ring\ chronology_t]$ , from 1929 to 2005 to maximize the  
247 overlap between the two sets of variables.

248 2. *Paleo Period* (1473 to 1929): Use  $f(\ln(Q_t) | PDSI_t)$  estimated in the previous step with the LBDA  
249 reconstructed PDSI,  $E[PDSI_t | tree-ring\ chronology_t]$ , to estimate  $\ln(Q_t)$  prior to 1929 (estimates of  $\ln(Q_t)$   
250 post 1929 will also be shown for comparison).

251 Several reconstruction model designs were considered given the initial diagnostic results. These  
252 included:

- 253 1) Developing one model for each streamflow station with predictors comprised of either:
  - 254 a. the LBDA grid located closest to the streamflow gauge;
  - 255 b. PCs of the LBDA from a region around the Missouri River Basin (Figure 1, green diamonds)  
256 within a multilinear model framework;
  - 257 c. canonical correlation analysis with regularization (rCCA described in Section 5.2) using LBDA  
258 grids within a 450km radius; or
  - 259 d. rCCA using LBDA grids within a 450km radius in addition to the retained PCs of CONUS-wide  
260 LBDA.

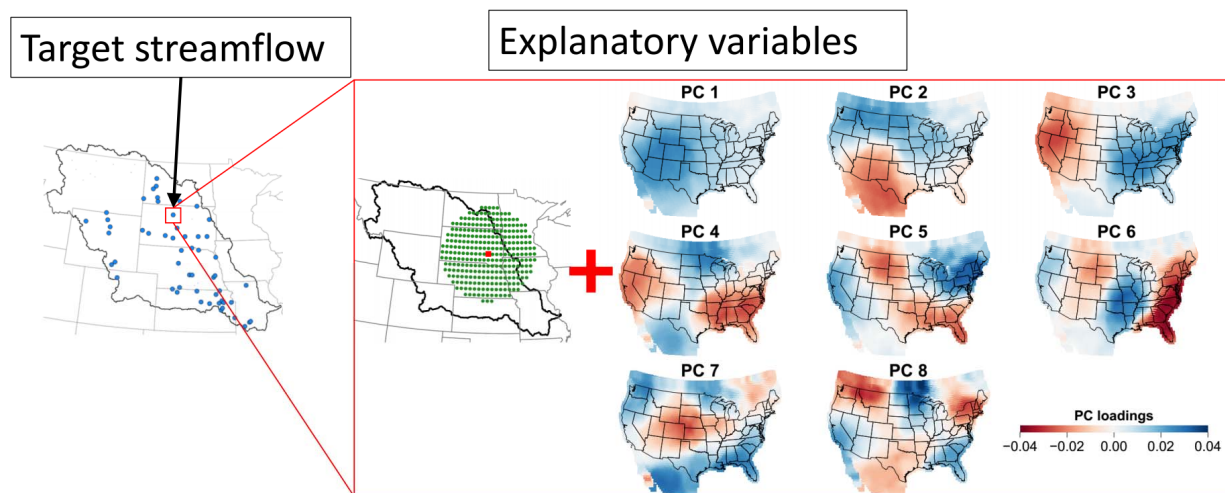
- 261 2) Developing one model for all streamflow stations with predictors comprised of either:  
262 a. PCs of both streamflow and CONUS-wide LBDA; or  
263 b. rCCA using streamflow and LBDA information from either the Missouri River Basin region or the  
264 CONUS region.

265 The fine resolution, gridded LBDA data results in a large number of highly correlated, potential  
266 predictors in all of the model designs considered with the exception of 1a. Given the  $0.5^\circ \times 0.5^\circ$   
267 resolution of the LBDA data, if a reconstruction of streamflow at a single gauge is considered using only  
268 a 450 km radius of surrounding LBDA data, one has around 300 potential predictors including the  
269 leading eight CONUS LBDA PCs. The length of annual streamflow records in the Missouri River Basin  
270 ranges from 39 to 85 years, and thus it is clear that the number of potential predictors is significantly  
271 greater than the number of observations available to fit the model. Of course, if one were to consider a  
272 model with a simultaneous reconstruction of all the streamflow records, the dimension of the  
273 estimation problem becomes even more challenging. Consequently, this is the first issue considered in  
274 model development.

275 A single reconstruction model would be advantageous in its ability to consider all streamflow stations at  
276 once rather than fitting 55 individual models. However, the varying streamflow record lengths and  
277 differences in record periods would have resulted in large uncertainties as a large degree of annual  
278 streamflow imputation would have been required. As a result, model designs under option 2 were not  
279 further examined.

280 A cursory comparison of the viable models was conducted by comparing the coefficient of  
281 determination for each fitted model. The development of individual models for each streamflow gauge  
282 using either model 1a (the closest LBDA grid) or 1b (PCs of the LBDA in a region around the Missouri  
283 River Basin) resulted in acceptable models of streamflow (mean adjusted  $R^2$  across the two models were

284 0.37 and 0.33 respectively). However, the spatially complete LBDA presents the opportunity to include a  
 285 wider variety of local and large-scale information and these models (model designs 1c and 1d) were  
 286 superior to the simpler models (model designs 1a and 1b). In addition, rCCA using both local and  
 287 CONUS-wide information (model 1d) resulted in slightly improved model results (mean adjusted  $R^2$  of  
 288 0.74 across all individual station models) over using only local information (model 1c, mean adjusted  $R^2$   
 289 of 0.71). We therefore selected model 1d (rCCA using LBDA grids within a 450km radius in addition to  
 290 the retained PCs of CONUS-wide LBDA) for further analysis. A schematic of the data used to fit this  
 291 model is shown in Figure 3, while a more detailed description of rCCA is provided in Section 5.2.



292  
 293 *Figure 3. Schematic of LBDA information to be included in the reconstruction model of one streamflow station in the Missouri*  
 294 *River Basin. The model inputs are the LBDA within a 450km radius and the first 8 PCs of US-wide LBDA.*

## 295 5.2 Regularized canonical correlation analysis

296 Methods such as principal component analysis (PCA, Jolliffe, 2002), archetype analysis (AA, Cutler and  
 297 Breiman, 1994; Stone and Cutler, 1996; Steinschneider and Lall, 2015) and canonical correlation analysis  
 298 (CCA, Hotelling, 1936) are typically used for dimension reduction in this setting. Given that the number  
 299 of potential predictors exceeds the number of observations, their high mutual correlation (>0.95 for  
 300 many of the adjacent grids), and an interest in exploring a multivariate streamflow response, we explore  
 301 the use of regularized canonical correlation where the regularization procedure is akin to ridge

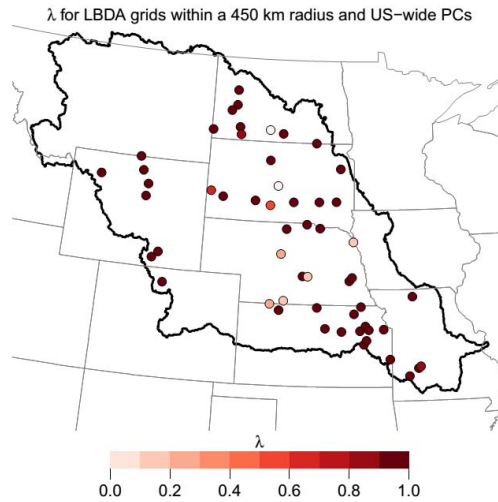
302 regression (De Bie and De Moor, 2003). Regularized canonical correlation was selected for use here to  
 303 capitalize on its ability to tailor the dimension reduction process to maximize the correlation between  
 304 the explanatory and target variables (in contrast to PCA where the explanatory and target variable  
 305 relationship is not considered).

306 Canonical correlation was introduced by Hotelling (1936) as a method of linearly transforming two  
 307 vector variables to canonical form to maximize the correlation between them. Consider two sets of  
 308 random variables represented by two matrices  $\mathbf{X}$  and  $\mathbf{Y}$ .  $\mathbf{X}$  is a  $n \times p$  matrix where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_p]$  containing  
 309  $n$  observations at  $p$  different locations, while  $\mathbf{Y}$  is a  $n \times q$  matrix where  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_q]$  containing  $n$   
 310 observations at  $q$  different locations. Both  $\mathbf{X}$  and  $\mathbf{Y}$  have finite variance matrices represented by  $\Sigma_{XX}$  and  
 311  $\Sigma_{YY}$  respectively. The covariance matrix between  $\mathbf{X}$  and  $\mathbf{Y}$  is  $\Sigma_{XY}$  while the covariance matrix between  $\mathbf{Y}$   
 312 and  $\mathbf{X}$  is  $\Sigma_{YX}$ . In this application,  $\mathbf{X}$  consists of the LBDA inputs specific to each streamflow gauge and  $\mathbf{Y}$  is  
 313 the streamflow at one station (i.e.  $p$  is large and  $q = 1$ ). Canonical correlation analysis involves rotating  
 314 the coordinate axes of both  $\mathbf{X}$  and  $\mathbf{Y}$  to new coordinate systems in order to clearly exhibit correlation  
 315 between  $\mathbf{X}$  and  $\mathbf{Y}$ . An arbitrary linear combination could be  $\mathbf{U} = \mathbf{X}\alpha$  and  $\mathbf{V} = \mathbf{Y}\gamma$  such that the correlation  
 316 between  $\mathbf{U}$  and  $\mathbf{V}$  is maximized.  $\mathbf{U}$  and  $\mathbf{V}$  yield the first pair of canonical variates, while  $\alpha$  and  $\gamma$  are the  
 317 vectors of canonical weights of length  $p$  and  $q$  respectively. This process may be repeated subject to the  
 318 constraint that following pairs of canonical variates are orthogonal to previous pairs with a maximum of  
 319  $\min(p, q)$  pairs obtained. In this application, the model of streamflow is applied station by station and  
 320 therefore only one pair of canonical variates are calculated and the first canonical variate of the LBDA is  
 321 used to fit the model of streamflow.

322 The correlation of successive pairs of canonical variates can be found using an eigen decomposition of  
 323  $\Sigma_{XX}^{-1}\Sigma_{XY}\Sigma_{YY}^{-1}\Sigma_{YX}$  and  $\Sigma_{YY}^{-1}\Sigma_{YX}\Sigma_{XX}^{-1}\Sigma_{XY}$ . The resulting  $\min(p, q)$  eigenvalues are common to both and the  
 324 square root of the eigenvalues yield the canonical correlation. The eigenvectors of  $\Sigma_{XX}^{-1}\Sigma_{XY}\Sigma_{YY}^{-1}\Sigma_{YX}$  and  
 325  $\Sigma_{YY}^{-1}\Sigma_{YX}\Sigma_{XX}^{-1}\Sigma_{XY}$  respectively yield the canonical weights  $\alpha$  and  $\gamma$  that are used to transform  $\mathbf{X}$  and  $\mathbf{Y}$ .

326 These weights are akin to the beta values in a multiple linear regression. Transforms using the  $i^{th}$   
327 eigenvector result in correlations corresponding to the square root of the  $i^{th}$  eigenvalue. Here, the  
328 weight for  $Y$  (streamflow) is 1 and the weights for  $X$  (LBDA) are given by the first eigenvector of  $\Sigma_{YY}^{-1}$   
329  $\Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}$ .

330 We employ regularization of the CCA process to address the issue of a large number of predictors  
331 relative to the number of observations (Vinod, 1976). Regularization is a smoothing process where a  
332 “roughness penalty”, also known as the regularization parameter ( $\lambda$ ), is introduced by converting  $\Sigma_{YY}$  and  
333  $\Sigma_{XX}$  to  $\Sigma_{YY} + \lambda_y I$  and  $\Sigma_{XX} + \lambda_x I$  respectively (Leurgans et al., 1993) and is similar to the technique of ridge  
334 regression (De Bie and De Moor, 2003). Values of  $\lambda$  range between 0 and 1 with larger  $\lambda$  values  
335 indicating a higher degree of smoothing. Regularization also enables the process of matrix inversion to  
336 be stabilized. A suitable value of  $\lambda$  is determined using a leave one out cross validation score, where  $\lambda_x$   
337 and  $\lambda_y$  are selected such that the correlation between the transformed datasets are maximized while the  
338 degrees of freedom used (as defined by Dijkstra, 2014) are limited to a maximum of  $n-10$ . If the criteria  
339 for the maximum degrees of freedom could not be achieved  $\lambda$  was set to one, otherwise  $\lambda$  was  
340 evaluated to two significant figures. No regularization was required for the single variable streamflow,  
341 whilst the LBDA was heavily regularized in almost all cases (Figure 4). rCCA was executed in R using the R  
342 package ‘CCA’ by González et al. (2008) and is freely available from the Comprehensive R Archive  
343 Network (CRAN <https://cran.r-project.org/>).



344

345 *Figure 4. CCA regularization parameter values for the explanatory variables (LBDA grids within a 450km radius and CONUS-wide*  
 346 *PCs) for each model of LBDA and streamflow.*

### 347 5.3 Model performance metrics

348 The performance of the model selected for further analysis includes verification using a leave-k-out  
 349 cross-validation procedure. Ten percent of the data (between 3 and 8 years out of a total of 32 and 77  
 350 years of streamflow data overlapping the LBDA record) are randomly selected and withheld from model  
 351 fitting. These values are then predicted from the model fit to the balance of the data. The entire process  
 352 is repeated 100 times, thus providing a set of 100 k-fold cross-validation samples. A comparison of the  
 353 model residuals resulting from both calibrated and verified model inputs is made for the 100 cross-  
 354 validation samples. The coefficient of efficiency (CE) and the reduction of error (RE) (Cook et al., 1994;  
 355 Wilson et al., 2010) are additional metrics calculated to verify the model. Both CE and RE are similar to  
 356 the Nash-Sutcliff efficiency test, however, the metrics are normalized using the mean of the verification  
 357 and calibration data respectively. Namely, CE is defined as

358

$$CE = 1 - \frac{\sum_{i=1}^n (x_i - \hat{x}_i)^2}{\sum_{i=1}^n (x_i - \bar{x}_v)^2}$$

359 while RE is defined as

$$RE = 1 - \frac{\sum_{i=1}^n (x_i - \hat{x}_i)}{\sum_{i=1}^n (x_i - \bar{x}_c)}$$

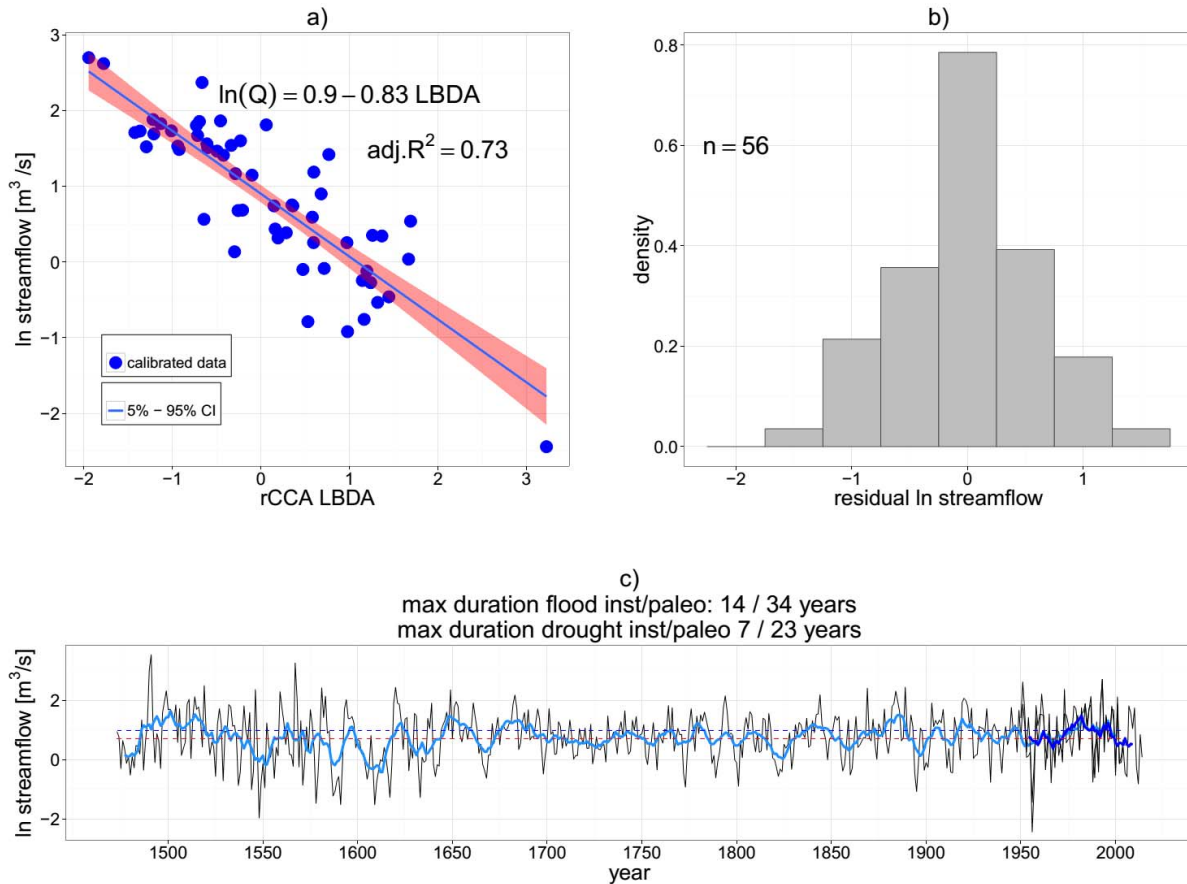
360 Where  $x_i$  and  $\hat{x}_i$  are the observed and modelled streamflows respectively, while  $\bar{x}_v$  and  $\bar{x}_c$  are the  
361 means of the observed streamflow in the validation and calibration periods respectively.

## 362 6 Results

363 A log-linear model was individually fitted to each streamflow gauge using logarithmically transformed  
364 streamflow and the first canonical variate of the LBDA inputs (schematic of inputs shown in Figure 3)  
365 using rCCA. The models were tested using a cross-validation procedure and calibration and verification  
366 metrics, as described in Section 5.3 were calculated. Finally, an overall assessment of the dominant  
367 modes of temporal variability in reconstructed streamflow variability in the Missouri River Basin was  
368 made using a frequency wavelet analysis on the leading PCs of reconstructed streamflow.

### 369 6.1 Model results

370 Streamflow at each of the 55 stations in the Missouri River Basin was constructed using a least squares  
371 regression model of natural-log streamflow and the first regularized canonical variate of the combined  
372 LBDA grids within a 450 km radius and first eight PCs of US-wide LBDA as the predictor variable. All  
373 results for modeled and reconstructed streamflow are shown in log space, while verification statistics  
374 are also calculated for logarithmic streamflow. An example of the input and modeled streamflow is  
375 shown for one streamflow station calibrated using data from Turkey Creek near Seneca (Figure 5). A  
376 summary of modeled streamflow results and summary statistics of the residuals for both calibrated and  
377 verified streamflow across all 100 calibration and verification sets are provided in the supporting  
378 information.



379

380 *Figure 5. Model results for Turkey Creek near Seneca (USGS gauge number 06814000) calibrated using all available data and the*  
 381 *first canonical variate of the LBDA inputs a) input (dots) and modelled (line) natural log streamflow showing 5<sup>th</sup> – 95<sup>th</sup> prediction*  
 382 *interval, b) histogram of model residuals, and c) flood and drought persistence gauged by a threshold of  $\text{mean} \pm 0.5SD$  of ten year*  
 383 *running average (blue lines – light blue line is the reconstructed ten year moving average).*

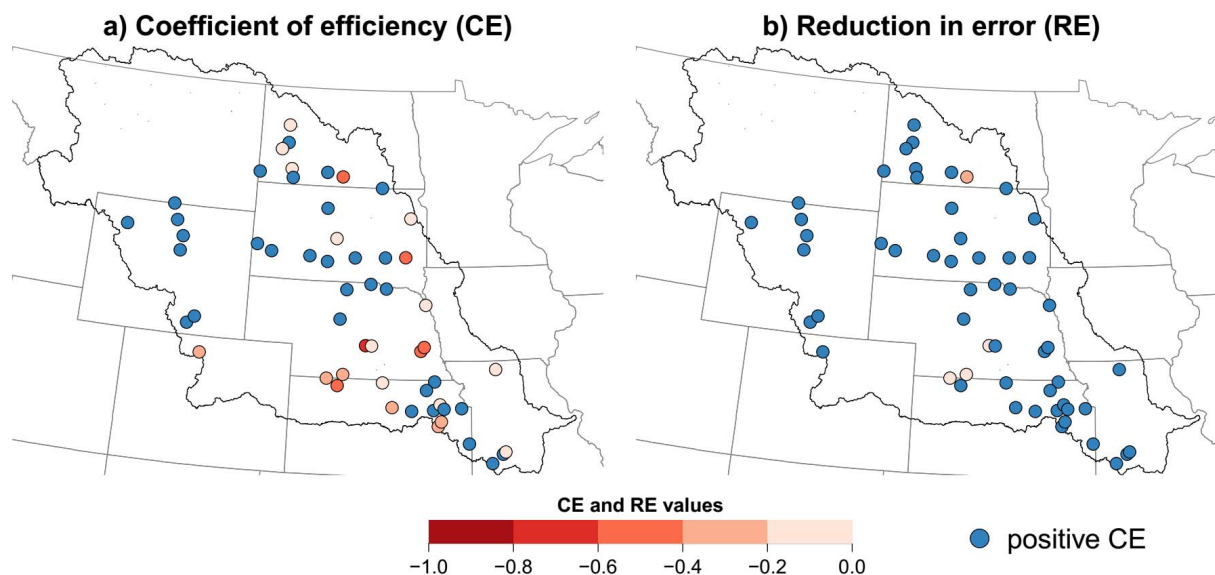
384 The rCCA linear model resulted in 55 models with adjusted  $R^2$  values ranging between 0.56 and 0.90.

385 The model residuals were near normal at almost all stations with the exception of some stations with  
 386 small catchment regions showing near-uniform residuals. In each streamflow model, the magnitude of  
 387 rCCA loadings for the eight PCs of CONUS LBDA were similar to those of the ~300 local LBDA grids  
 388 indicating that information from the eight PCs of CONUS LBDA did not dominate the reconstruction.

## 389 6.2 Model validation

390 The predictive power of the model was assessed by calculating the coefficient of efficiency (CE) and  
 391 reduction of error (RE) for all cross validations. The median CE and RE values are shown in Figure 6 a and

392 b respectively, while distributions of the values are shown in box plots in the supporting information. CE  
393 and RE values range from  $-\infty$  to 1, with values over 0 indicating that the model predictions are more  
394 accurate than the respective climatology used for each statistic.

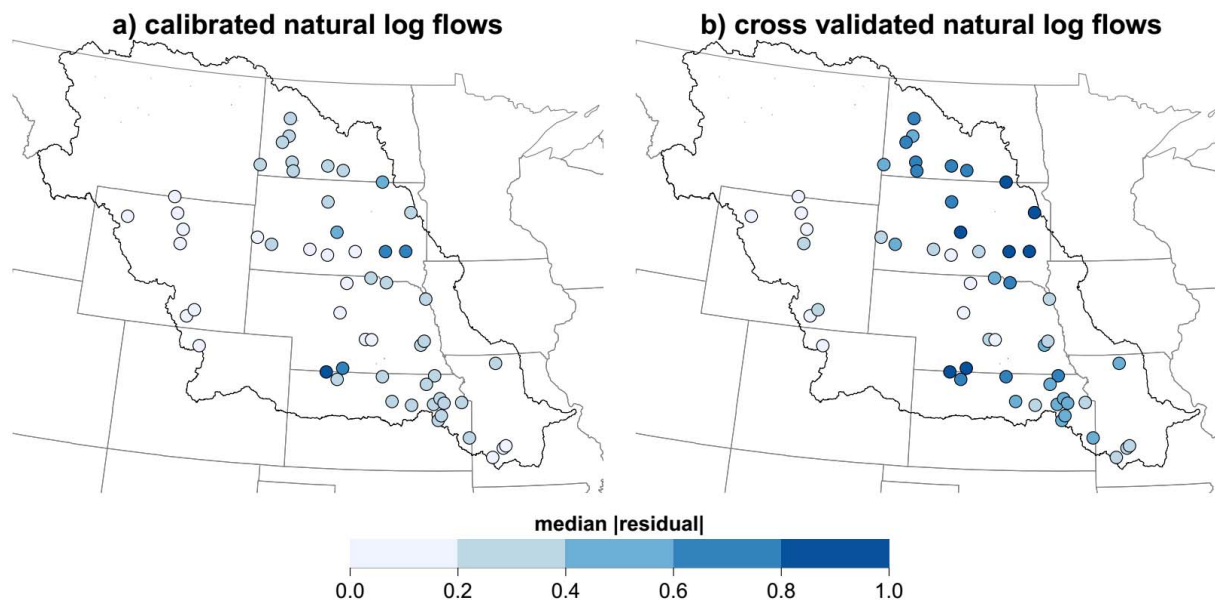


395  
396 *Figure 6. Median values of a) CE and b) RE values for cross validated models of the 55 Missouri River Basin streamflow stations.*

397 The CE values in Figure 6 indicate that many of the streamflow models are able to provide median  
398 streamflow predictions with greater accuracy than the climatology of the withheld data. The RE values  
399 likewise suggest that the models are able to predict the withheld values with greater accuracy than the  
400 climatology of the calibrated values at most stations. The combined CE and RE results suggest there is  
401 some skill in reconstructing paleoclimate streamflow using the LBDA.

402 A comparison of the distribution of residuals resulting from calibrated values vs. validated values was  
403 also made using the 100 repetitions of the k-fold cross validation test. The median and the range of the  
404 cross-validated absolute residuals are modestly larger than that of the absolute residuals during the  
405 fitting process. The differences, accounting for the sample sizes, are not statistically significant in over  
406 90% of 100 cross validated results (using a two-sided t-test and a null hypothesis that the true difference  
407 between the means is zero and  $\alpha = 0.05$ ). The median values of the cross validated residuals are

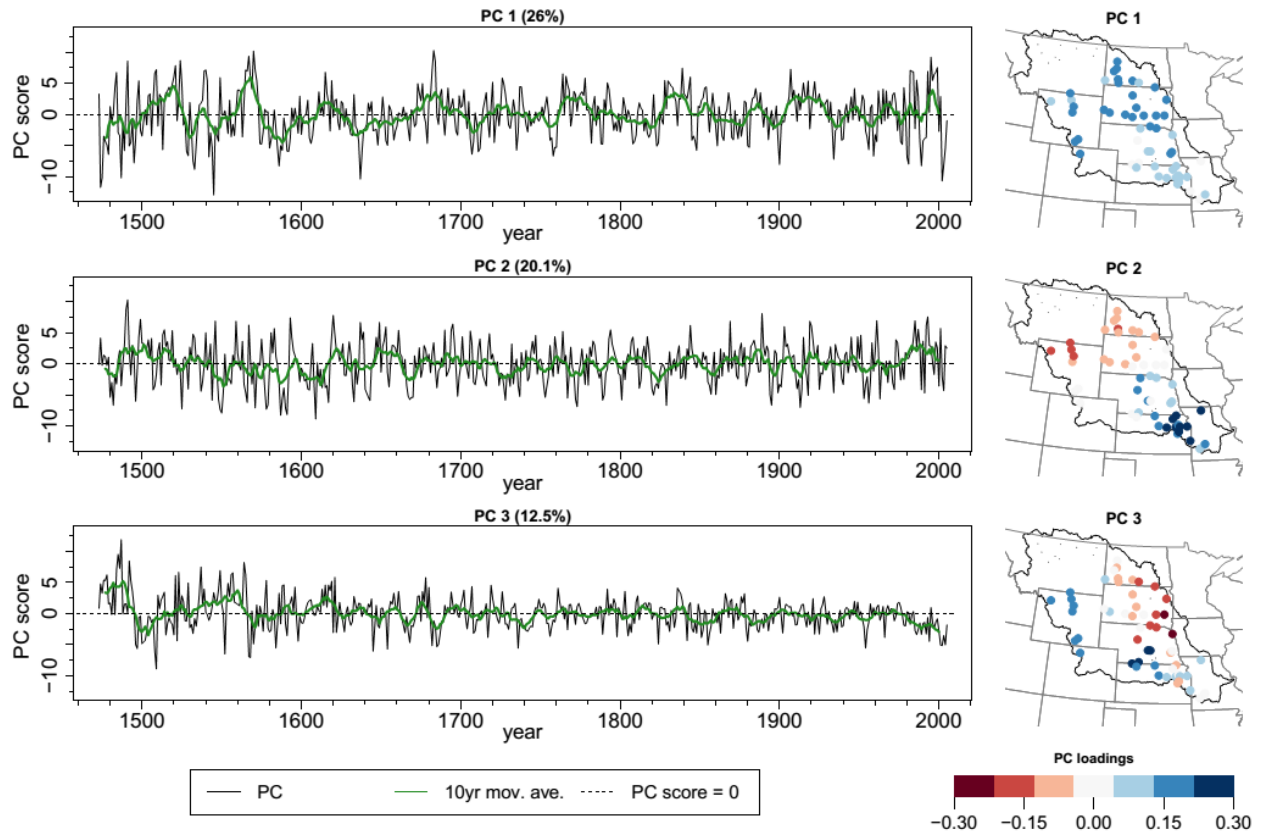
408 shown in Figure 7, while box plots showing the distribution of the cross validated residuals are shown in  
409 the supporting information. A summary of the reconstruction is presented in Section 6.3 using PCs of the  
410 streamflow reconstructed at the 55 stations.



411  
412 *Figure 7. Median cross validated model residuals (absolute value) for a) calibrated and b) verified natural log streamflow inputs*  
413 *at the 55 Missouri River Basin streamflow stations. Cross validations were repeated 100 times.*

### 414 6.3 Reconstructed streamflow analysis

415 Mean annual streamflow was reconstructed for all 55 stations in the Missouri River Basin (data available  
416 in Ho et al., 2016). PCA was used to extract the leading modes of variability from the reconstructed  
417 Missouri River Basin natural log streamflow at all 55 stations. A combined assessment of the PCs using  
418 North's Rule of Thumb and a scree plot showed adjacent degenerate eigenvalues pairs for PCs of order 4  
419 and higher and discontinuities at PC4 respectively (shown in supporting information). Furthermore, the  
420 spatial pattern of PC4 was difficult to interpret and therefore only the first three PCs are shown in Figure  
421 8.



422

423 *Figure 8. Annual (black) and 10 year moving average (green line) of the first four PCs of reconstructed mean natural log*  
 424 *streamflow in the Missouri River Basin (left) with percentage variance explained shown in parenthesis and the corresponding*  
 425 *loading patterns (right).*

426 Negative streamflow anomalies coinciding with the 1930s dust bowl drought and 1950s drought are  
 427 evident in PC1, which has a positive loading pattern across all stations in the Missouri River Basin (Figure  
 428 8). The 1950s drought had more severe impacts in the southern half of the Missouri River Basin  
 429 (Piechota and Dracup, 1996; Cole et al., 2002; Andreadis et al., 2005; Cook et al., 2009) and this is  
 430 detected in PC2 that is comprised of positive loadings in the southern half of the Basin. The Civil War  
 431 drought, which largely impacted the Central Plains region from the mid-1850s to mid-1860s (Herweijer  
 432 et al., 2006), is also evident in PC2.

433 The severity of these droughts had wide ranging impacts on ecological states, agricultural produce and  
 434 social activities and are often used as benchmark droughts to which contemporary droughts are  
 435 compared (Breshears et al., 2005; Hornbeck, 2009). However, the streamflow reconstructions suggest

436 periods where streamflow deficits may have been more severe than any of these historical droughts. For  
437 example, the late 1540s, 1590s, and late 1750s all show negative streamflow anomalies in both PC1 and  
438 PC2 (26% and 20.1% variability explained respectively) that are of similar or greater magnitude than the  
439 streamflow deficits associated with the Civil War, Dust Bowl or 1950s droughts. Drought durations were  
440 also longer when assessed in the context of streamflow variability over the past 500 years. A threshold  
441 of above or below 0.5 standard deviation in decadal streamflow at each station was used as an  
442 approximation of flood and drought regimes respectively. The maximum duration of continuous periods  
443 of either flood or drought regimes were found to be longer in the majority of stations investigated as  
444 demonstrated in Figure 5 c that shows the reconstructed annual and decadal streamflow at Turkey  
445 Creek and the longest duration drought and flood regime for the station.

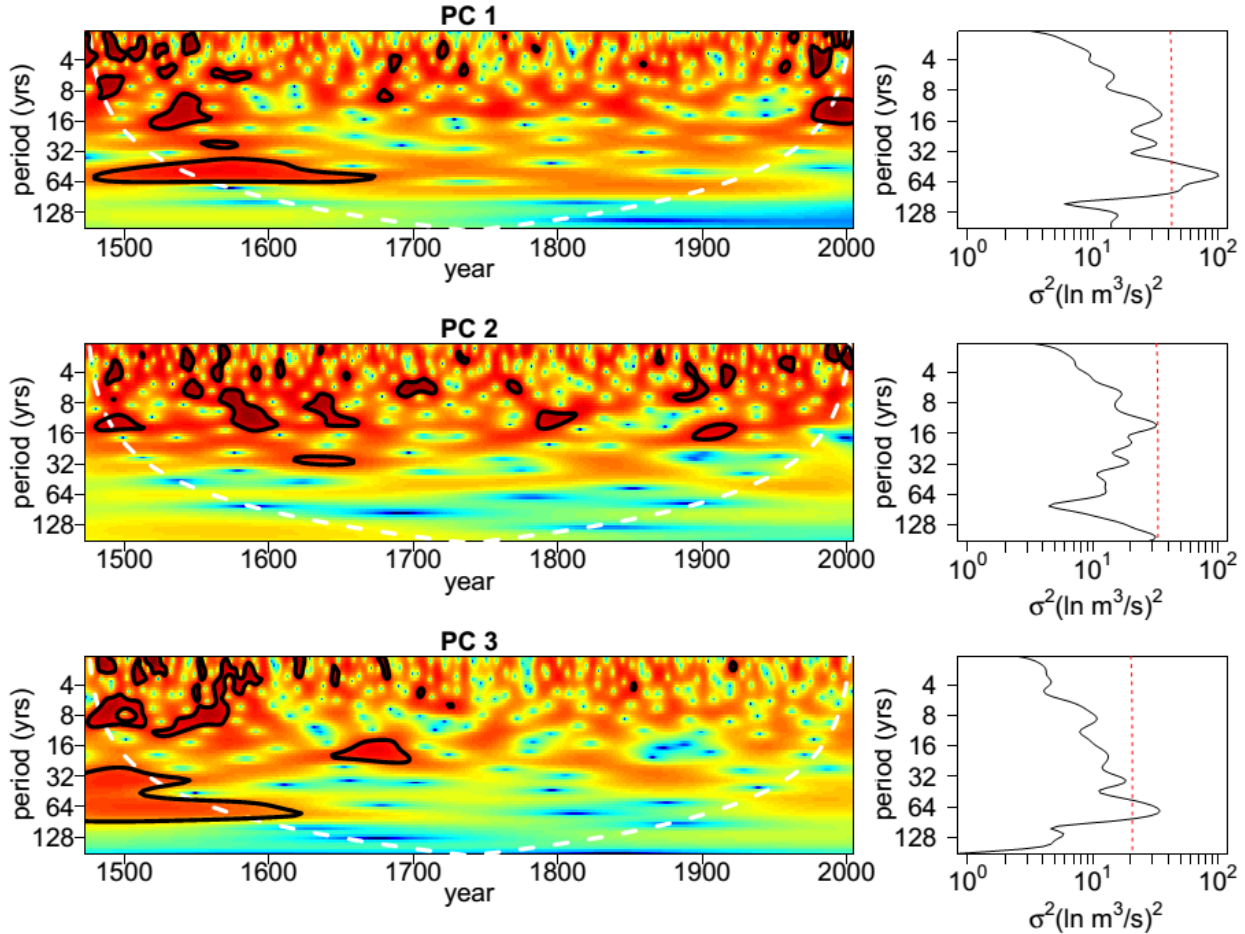
446 PC3 (12.5% variability explained) is a mode of variability largely influenced by differences in streamflow  
447 variability between the Rocky Mountains regions and the remainder of the Missouri River Basin. The  
448 apparent increase in variability in the 1500s and 1600s in PC3 suggests that differences in streamflow  
449 variability between the Rocky Mountains regions and the plains region may have been more  
450 pronounced during these two centuries.

451 PC2 of the reconstructed streamflow has positive loadings in the southern Missouri River Basin. This PC  
452 shows an extended period of below average flows, with the exception of a couple of years, in the first  
453 two decades of the 1600s and is evident in individual streamflow reconstructions around the border of  
454 Nebraska and South Dakota and in Kansas. A similar persistence of decadal-scale low flows are not  
455 featured in the instrumental record, suggesting that successive years of low streamflow have persisted  
456 for longer periods than what has been observed in historical records.

457 The PCs of the reconstructed streamflow also show several periods where annual flows are larger than  
458 the maximum instrumental streamflow. Periods of anomalously high streamflow include the 1560s for

459 PC 1 and the early 1490s for PC 2 and 3. These periods coincide with high positive reconstructed PDSI  
460 values particularly in the lower Missouri Basin region and along the middle and southern Rocky  
461 Mountains in the early 1490s and positive PDSI values across most of the interior plains. Anomalously  
462 large streamflow during the 1560s were reconstructed at stream gauges around the interior plains  
463 regions from tributaries that join the Missouri River in North and South Dakota (e.g. Cannonball, Grand  
464 and White Rivers). However, anomalously high streamflow around the 1560s was not reconstructed for  
465 either the Platte or Yellowstone rivers suggesting that the increase in streamflow may have resulted  
466 from weather systems delivering moisture in the upper Interior Plains (North and South Dakota), but not  
467 in the Rocky Mountains region.

468 In order to quantitatively assess the temporal modes of variability implied by the PCs, a continuous  
469 wavelet transform was applied to identify key periodicities in the main modes of reconstructed  
470 streamflow variability. Wavelet transforms for PCs 1-3 along with their global wavelet spectra are shown  
471 in Figure 9. The wavelet transforms were implemented using the Morlet wavelet function as the mother  
472 wavelet and padded with zeros to limit the edge effects of the wavelet analysis (Torrence and Compo,  
473 1998).



475

476 *Figure 9. Wavelet transforms of PC 1-3 of reconstructed log streamflow in the Missouri River Basin from 1473-2005 (left), black*  
 477 *lines show 95% significance level, dashed white line shows cone of influence under which edge effects may influence the results*  
 478 *and the time averaged global wavelet spectra (right) with red dashed line showing 95% confidence level for the wavelet*  
 479 *transform using a white noise background spectrum.*

480 All PCs show decadal-scale variability in the 1500s and 1600s, which is significant at the 95% level  
 481 against white noise for PC 1 and 3. Interestingly, none of the wavelet spectra show a consistent mode of  
 482 variability that coincide with an El Niño Southern Oscillation frequency (2-7 years Allan, 2000). In  
 483 addition, previously identified periods of decadal-scale variability in paleoclimate reconstructions of the  
 484 Pacific Decadal Oscillation (Gedalof and Smith, 2001; MacDonald and Case, 2005) are not reproduced in  
 485 the streamflow reconstruction. The multi-decadal-scale variability seen in PC1 and PC3 around the 16<sup>th</sup>  
 486 century suggest that the persistence of drought and pluvial events was longer around this period  
 487 compared with historic records and is reflected in previous reconstructions of persistent drought during

488 this period (Woodhouse and Overpeck, 1998) and in nearby regions (Meko et al., 1995). A similar  
489 pattern of multi-decadal-scale variability is also seen in the PCs of LBDA grids within the Missouri River  
490 Basin around the 16<sup>th</sup> century (results not shown here).

## 491 7 Discussion and Conclusions

492 Human activities and social processes are irrefutably intertwined with the hydrological cycle with each  
493 influencing and affecting the behavior of the other (Matalas et al., 1982). Consequently, hydrologic  
494 investigations need to occur at spatial scales that encompass the social and political regions influencing  
495 how water is stored, used, and transported. This requirement necessitates a broader spatial scale to be  
496 considered beyond the watershed scale typically addressed in hydrologic assessments (Vogel et al.,  
497 2015). Spatially and temporally broad measures of hydroclimatic variability such as the LBDA provide  
498 opportunities to conduct such assessments.

499 We used the LBDA to demonstrate a method of reconstructing mean annual streamflow in the Missouri  
500 River Basin from 1473-2005 in order to facilitate a future reconstruction of streamflow across the  
501 CONUS. The rCCA approach enabled a large degree of spatially coherent information in the LBDA to be  
502 condensed via linear transforms into a canonical variate that explained the majority of information in  
503 the target streamflow. The additional regularization step addressed the issue of an ill-posed problem  
504 where the number of variables far exceeded the number of observations. Regularization addressed  
505 potential issues with overfitting of the model and subsequent poor predictability of uncalibrated values.  
506 This site-by-site approach, opposed to a multi-site model, avoided data issues where streamflow records  
507 were not concurrently available at all sites. The method also tailored the LBDA information included in  
508 the model to each unique site and therefore has the potential to be applied to any individual streamflow  
509 site in the United States and, potentially, any streamflow site in North America where there are useful  
510 reconstructions in the LBDA.

511 The analysis of streamflow variability in the Missouri River Basin over the past 500 years was, however,  
512 limited by the absence of streamflow gauges in the headwaters of the Missouri River. Based on the  
513 selection criteria used here, there are no gauges located upstream of Fort Peck Lake and only one gauge  
514 in Montana was included in the analysis. The dearth of gauges in this region therefore omits potentially  
515 critical information when drawing conclusions regarding regional streamflow variability as presented in  
516 Section 6.3. Future analyses could seek to include additional streamflow gauges, potentially drawing on  
517 non-references gauges from the GAGESII database and other sources. While non-reference GAGESII  
518 gauges are located within disturbed catchments with potentially regulated streams, it may be plausible  
519 that these measurements would still preserve the signatures of interannual streamflow variability which  
520 we seek to recover. Encouragingly, the headwater regions of the Missouri River Basin contain a  
521 relatively dense network of tree-ring chronologies that would improve the fidelity of both PDSI and  
522 subsequent streamflow reconstructions. Furthermore, efforts towards expanding this tree-ring network  
523 are currently being made with a view to target and improve Missouri River Basin streamflow  
524 reconstructions (Pederson, 2013).

525 Dendrochronological studies for the purposes of climate reconstruction typically use a targeted  
526 sampling method whereby old living trees are selected for sampling to maximize the length of  
527 reconstruction and the location of sampling is chosen to ensure that tree growth is largely limited by the  
528 climatic factors of interest for reconstruction (Speer, 2010). Over 50 different tree species were used in  
529 the formulation of the LBDA, however, species was not a determinant for inclusion in the LBDA. Rather,  
530 the tree-ring chronologies specifically selected for use in reconstruction were those that were best  
531 correlated with available soil moisture as modeled a priori by PDSI and this was done using a method of  
532 point-by-point regression that ensures that the selected chronologies were likely to have a causal  
533 association with the gridded instrumental PDSI being reconstructed (Cook et al., 1999).

534 In developing our approach to reconstructing streamflow from the tree-ring based LBDA, we have  
535 capitalized on the fact that both variables are derivatives of a set of (unspecified) climate variables to  
536 which both streamflow and PDSI are sensitive (e.g. precipitation amount, temperature, and  
537 evaporation). Although variability is lost in each step of reconstructing PDSI from tree-rings and in  
538 reconstructing streamflow from the LBDA, it is possible that the reconstructed LBDA may remove noise  
539 irrelevant to hydrological variability rendering this paleoclimate reconstruction of PDSI particularly  
540 suitable for informing streamflow. Furthermore, streamflow and PDSI are inherently different  
541 hydrological variables as demonstrated in the temporal and spatial statistics shown in the supporting  
542 information. Streamflow is also a derivative of a suite of non-climatic variables such as catchment size,  
543 land-use, geology, topography, and groundwater interactions. These catchment-specific attributes are  
544 not explicitly considered in our study (with the exception of land use considered through the selection of  
545 streamflow gauges from relatively undisturbed catchments) and the inclusion of attributes could  
546 provide an improvement to the model (e.g. Thomas and Benson, 1970; Lima and Lall, 2010).

547 Assessments of hydrological variability on temporal scales broader than traditional hydrological  
548 practices are paramount to identifying long-term variations and quantifying the nature of persistent  
549 regimes, which are unachievable using comparatively short instrumental record. The understanding of  
550 paleoclimate variability facilitates the development of water and land use practices that are appropriate,  
551 sustainable, and resilient under previous patterns of hydroclimatic variability and complement efforts  
552 towards developing suitable adaption procedures for projected future climate scenarios. Whilst  
553 paleoclimate reconstructions of annual streamflow are useful for assessing long-term variability, annual  
554 information is often too coarse for water resource applications such as reservoir management and  
555 seasonal water allocations. In order to produce relevant data on a monthly or daily time step, a future  
556 undertaking could involve temporal disaggregation of the data. One method that could be utilized is a  
557 non-parametric k-nearest neighbors approach (Lall and Sharma, 1996; Rajagopalan and Lall, 1999;

558 Nowak et al., 2010) that would avoid some issues associated with multi-site stochastic generation  
559 schemes (Valencia and Schaake, 1973; Segond et al., 2006).The objective of producing paleoclimate  
560 reconstructions of streamflow that are relevant for water resource management also suggests that the  
561 reconstruction of streamflow that is currently regulated would be valuable. Paleoclimate  
562 reconstructions of regulated streamflow could be achieved at an annual scale provided that regulations  
563 primarily impact the timing of sub-annual flows or if sufficient data prior to regularization exists. The  
564 resulting annual streamflow reconstruction could be disaggregated to a shorter timescale if pre-  
565 regulated flow data is available.

566 This study has capitalized on the availability of the spatially and temporally complete LBDA dataset to  
567 attempt a reconstruction of streamflow on a broad spatial scale across the Missouri River Basin. The  
568 analysis employed here provides a feasible foundation from which streamflow reconstructions across  
569 the CONUS and North America could be implemented using the LBDA. In addition, the methodology  
570 presented here could be applied to reconstructing streamflow using other reconstructions of drought  
571 including the Monsoon Asia Drought Atlas (Cook et al., 2010b), the Old World Drought Atlas covering  
572 Europe, North Africa, and the Middle East (Cook et al., 2015b) and the Australia and New Zealand  
573 summer drought atlas (Palmer et al., 2015).

574

## 575 Acknowledgments

576 We would like to sincerely thank Ben Cook (bc9z@ldeo.columbia.edu) for providing the LBDA data, the  
577 US Geological Survey (USGS) for making their monthly streamflow data available  
578 ([http://waterdata.usgs.gov/nwis/monthly?referred\\_module=sw](http://waterdata.usgs.gov/nwis/monthly?referred_module=sw)), staff at the USGS (Christopher Ryan  
579 cmryan@usgs.gov and Thomas Weaver tlweaver@usgs.gov) for clarifying streamflow data availability;  
580 Siyan Wang for initial data collection; Scott Steinschneider, Xun Sun and Pierre Gentine at the Columbia  
581 Water Center for their invaluable discussions on this work; and Juan A. Ballesteros and one other  
582 anonymous reviewer for their constructive comments. Reconstruction results are available on the NOAA  
583 paleoclimate data base (<https://www.ncdc.noaa.gov/paleo/study/19520>). This work is funded by an NSF  
584 award 1360446 and NSF award 1401698. Lamont-Doherty contribution number XXXX.

## 585 References

- 586 Allan, R. J., 2000: ENSO and Climatic Variability in the Past 150 Years. *El Niño and the Southern*  
587 *Oscillation: Multiscale Variability and Global and Regional Impacts*, Henry F. Diaz and Vera  
588 Markgraf, Eds., Cambridge University Press, pp.
- 589 Allen, E. B., Rittenour, T. M., DeRose, R. J., Bekker, M. F., Kjelgren, R. and Buckley, B. M., 2013: A tree-  
590 ring based reconstruction of Logan River streamflow, northern Utah, *Water Resources Research*,  
591 **49** (12), 8579-8588, doi: 10.1002/2013wr014273.
- 592 Andreadis, K. M., Clark, E. A., Wood, A. W., Hamlet, A. F. and Lettenmaier, D. P., 2005: Twentieth-  
593 Century Drought in the Conterminous United States, *Journal of Hydrometeorology*, **6** (6), 985-  
594 1001, doi: 10.1175/JHM450.1.
- 595 Ault, T. R., Cole, J. E., Overpeck, J. T., Pederson, G. T. and Meko, D. M., 2014: Assessing the risk of  
596 persistent drought using climate model simulations and paleoclimate data, *Journal of Climate*,  
597 **27** (20), 7529-7549, doi: 10.1175/jcli-d-12-00282.1.
- 598 Breshears, D. D., Cobb, N. S., Rich, P. M., Price, K. P., Allen, C. D., Balice, R. G., Romme, W. H., Kastens, J.  
599 H., Floyd, M. L., Belnap, J., Anderson, J. J., Myers, O. B. and Meyer, C. W., 2005: Regional  
600 vegetation die-off in response to global-change-type drought, *Proceedings of the National*  
601 *Academy of Sciences of the United States of America*, **102** (42), 15144-15148, doi:  
602 10.1073/pnas.0505734102.
- 603 Brown, D. P. and Comrie, A. C., 2004: A winter precipitation 'dipole' in the western United States  
604 associated with multidecadal ENSO variability, *Geophysical Research Letters*, **31** (9), doi:  
605 10.1029/2003GL018726.
- 606 Buuren, S. v. and Groothuis-Oudshoorn, K., 2011: mice: Multivariate Imputation by Chained Equations in  
607 R, *Journal of Statistical Software*, **45** (3), 67.

608 Cole, J. E., Overpeck, J. T. and Cook, E. R., 2002: Multiyear La Niña events and persistent drought in the  
609 contiguous United States, *Geophysical Research Letters*, **29** (13), 25-1-25-4, doi:  
610 10.1029/2001GL013561.

611 Cook, B. I., Miller, R. L. and Seager, R., 2009: Amplification of the North American “Dust Bowl” drought  
612 through human-induced land degradation, *Proceedings of the National Academy of Sciences*,  
613 **106** (13), 4997-5001, doi: 10.1073/pnas.0810200106.

614 Cook, B. I., Ault, T. R. and Smerdon, J. E., 2015a: Unprecedented 21st century drought risk in the  
615 American Southwest and Central Plains, *Science Advances*, **1** (1), doi: 10.1126/sciadv.1400082.

616 Cook, B. I., Smerdon, J. E., Seager, R. and Cook, E. R., 2013a: Pan-Continental Droughts in North America  
617 over the Last Millennium, *Journal of Climate*, **27** (1), 383-397, doi: 10.1175/jcli-d-13-00100.1.

618 Cook, E. R. and Krusic, P. J., 2004: The North American Drought Atlas,  
619 Cook, E. R., Briffa, K. R. and Jones, P. D., 1994: Spatial regression methods in dendroclimatology: A  
620 review and comparison of two techniques, *International Journal of Climatology*, **14** (4), 379-402,  
621 doi: 10.1002/joc.3370140404.

622 Cook, E. R., Meko, D. M., Stahle, D. W. and Cleaveland, M. K., 1999: Drought Reconstructions for the  
623 Continental United States, *Journal of Climate*, **12** (4), 1145-1162, doi: 10.1175/1520-  
624 0442(1999)012<1145:drftcu>2.0.co;2.

625 Cook, E. R., Woodhouse, C. A., Eakin, C. M., Meko, D. M. and Stahle, D. W., 2004: Long-Term Aridity  
626 Changes in the Western United States, *Science*, **306** (5698), 1015-1018, doi:  
627 10.1126/science.1102586.

628 Cook, E. R., Seager, R., Heim, R. R., Vose, R. S., Herweijer, C. and Woodhouse, C., 2010a: Megadroughts  
629 in North America: placing IPCC projections of hydroclimatic change in a long-term palaeoclimate  
630 context, *Journal of Quaternary Science*, **25** (1), 48-61, doi: 10.1002/jqs.1303.

631 Cook, E. R., Anchukaitis, K. J., Buckley, B. M., D’Arrigo, R. D., Jacoby, G. C. and Wright, W. E., 2010b:  
632 Asian Monsoon Failure and Megadrought During the Last Millennium, *Science*, **328** (5977), 486-  
633 489, doi: 10.1126/science.1185188.

634 Cook, E. R., Palmer, J. G., Ahmed, M., Woodhouse, C. A., Fenwick, P., Zafar, M. U., Wahab, M. and Khan,  
635 N., 2013b: Five centuries of Upper Indus River flow from tree rings, *Journal of Hydrology*, **486**  
636 (0), 365-375, doi: 10.1016/j.jhydrol.2013.02.004.

637 Cook, E. R., Seager, R., Kushnir, Y., Briffa, K. R., Büntgen, U., Frank, D., Krusic, P. J., Tegel, W., van der  
638 Schrier, G., Andreu-Hayles, L., Baillie, M., Baittinger, C., Bleicher, N., Bonde, N., Brown, D.,  
639 Carrer, M., Cooper, R., Čufar, K., Dittmar, C., Esper, J., Griggs, C., Gunnarson, B., Günther, B.,  
640 Gutierrez, E., Haneca, K., Helama, S., Herzig, F., Heussner, K.-U., Hofmann, J., Janda, P., Kontic,  
641 R., Köse, N., Kyncl, T., Levanič, T., Linderholm, H., Manning, S., Melvin, T. M., Miles, D., Neuwirth,  
642 B., Nicolussi, K., Nola, P., Panayotov, M., Popa, I., Rothe, A., Seftigen, K., Seim, A., Svarva, H.,  
643 Svoboda, M., Thun, T., Timonen, M., Touchan, R., Trotsiuk, V., Trouet, V., Walder, F., Ważny, T.,  
644 Wilson, R. and Zang, C., 2015b: Old World megadroughts and pluvials during the Common Era,  
645 *Science Advances*, **1** (10), doi: 10.1126/sciadv.1500561.

646 Cutler, A. and Breiman, L., 1994: Archetypal Analysis, *Technometrics*, **36** (4), 338-347, doi:  
647 10.1080/00401706.1994.10485840.

648 Dai, A., 2013: Increasing drought under global warming in observations and models, *Nature Clim.*  
649 *Change*, **3** (1), 52-58, doi: 10.1038/nclimate1633.

650 De Bie, T. and De Moor, B., 2003: On the regularization of canonical correlation analysis, *Int. Sympos. ICA*  
651 *and BSS*, 785-790.

652 Dettinger, M. D., Cayan, D. R., Diaz, H. F. and Meko, D. M., 1998: North–South Precipitation Patterns in  
653 Western North America on Interannual-to-Decadal Timescales, *Journal of Climate*, **11** (12), 3095-  
654 3111, doi: 10.1175/1520-0442(1998)011<3095:NSPPIW>2.0.CO;2.

655 Devineni, N., Lall, U., Pederson, N. and Cook, E., 2013: A Tree-Ring-Based Reconstruction of Delaware  
656 River Basin Streamflow Using Hierarchical Bayesian Regression, *Journal of Climate*, **26** (12),  
657 4357-4374, doi: 10.1175/jcli-d-11-00675.1.

658 Dijkstra, T., 2014: Ridge regression and its degrees of freedom, *Quality & Quantity*, **48** (6), 3185-3193,  
659 doi: 10.1007/s11135-013-9949-7.

660 Driscoll, D., 2013: *Hydroclimatic information and risk assessment methods for managing hydrologic*  
661 *extremes in the Missouri River Basin - A white paper for future studies. Hydroclimatic extremes in*  
662 *the Missouri River Basin*, USGS.

663 Fritts, H. C., 1976: *Tree Rings and Climate*. Academic Press Inc., 567 pp.

664 Galat, D. L., Berry Jr, C. R., Peters, E. J. and White, R. G., 2005: Chapter 10 - Missouri River Basin. *Rivers of*  
665 *North America*, Arthur C. Benke and Colbert E. Cushing, Eds., Academic Press, pp. 426-480.

666 Gallant, A. J. E. and Gergis, J., 2011: An experimental streamflow reconstruction for the River Murray,  
667 Australia, 1783-1988, *Water Resources Research*, **47** (12), W00G04, doi:  
668 10.1029/2010wr009832.

669 Gedalof, Z. e. and Smith, D. J., 2001: Interdecadal climate variability and regime-scale shifts in Pacific  
670 North America, *Geophys. Res. Lett.*, **28** (8), 1515-1518, doi: 10.1029/2000gl011779.

671 Gershunov, A. and Barnett, T. P., 1998: Interdecadal Modulation of ENSO Teleconnections, *Bulletin of*  
672 *the American Meteorological Society*, **79** (12), 2715-2725, doi: 10.1175/1520-  
673 0477(1998)079<2715:IMOET>2.0.CO;2.

674 González, I., Déjean, S., Martin, P. G. P. and Baccini, A., 2008: CCA: An R Package to Extend Canonical  
675 Correlation Analysis, *2008*, **23** (12), 14, doi: 10.18637/jss.v023.i12.

676 Herweijer, C., Seager, R. and Cook, E. R., 2006: North American droughts of the mid to late nineteenth  
677 century: a history, simulation and implication for Mediaeval drought, *The Holocene*, **16** (2), 159-  
678 171, doi: 10.1191/0959683606hl917rp.

679 Herweijer, C., Seager, R., Cook, E. R. and Emile-Geay, J., 2007: North American Droughts of the Last  
680 Millennium from a Gridded Network of Tree-Ring Data, *Journal of Climate*, **20** (7), 1353-1376,  
681 doi: 10.1175/JCLI4042.1.

682 Higgins, R. W., Yao, Y., Yarosh, E. S., Janowiak, J. E. and Mo, K. C., 1997: Influence of the Great Plains  
683 Low-Level Jet on Summertime Precipitation and Moisture Transport over the Central United  
684 States, *Journal of Climate*, **10** (3), 481-507, doi: 10.1175/1520-  
685 0442(1997)010<0481:IOTGPL>2.0.CO;2.

686 Ho, M., Kiem, A. S. and Verdon, D. C., 2015: A paleoclimate rainfall reconstruction in the Murray-Darling  
687 Basin (MDB), Australia: 2. Assessing hydroclimatic risk using paleoclimate records of wet and dry  
688 epochs, *Water Resour. Res.*, **51** (10), doi: 10.1002/2015WR017059.

689 Ho, M., Lall, U. and Cook, E. R., 2016: Missouri River Basin 533 Year Annual Streamflow  
690 Reconstructions, <<https://www.ncdc.noaa.gov/paleo/study/19520>>.

691 Hornbeck, R., 2009: The Enduring Impact of the American Dust Bowl: Short and Long-run Adjustments to  
692 Environmental Catastrophe, *National Bureau of Economic Research Working Paper Series*, **No.**  
693 **15605**, doi: 10.3386/w15605.

694 Hotelling, H., 1936: Relations between two sets of variates, *Biometrika*, **28** (3-4), 321-377, doi:  
695 10.1093/biomet/28.3-4.321.

696 Jolliffe, I. T., 2002: *Principal Component Analysis*. Second edition ed. Springer-Verlag, 487 pp.

697 Jones, P. D. and Mann, M. E., 2004: Climate over past millennia, *Rev. Geophys.*, **42** (2), RG2002, doi:  
698 10.1029/2003rg000143.

699 Kunkel, K. E., Stevens, L. E., Stevens, S. E., Sun, L., Janssen, E., Wuebbles, D., Kruk, M. C., Thomas, D. P.,  
700 Shulski, M., Umphlett, N., Hubbard, K., Robbins, K., Romolo, L., Akyuz, A., Pathak, T., Bergantino,  
701 T. and Dobson, J. G., 2013: Regional Climate Trends and Scenarios for the U.S. National Climate  
702 Assessment. Part 4. Climate of the U.S. Great Plains, NOAA Technical Report, NOAA, pp. 82.

703 Lall, U. and Sharma, A., 1996: A Nearest Neighbor Bootstrap For Resampling Hydrologic Time Series,  
704 *Water Resources Research*, **32** (3), 679-693, doi: 10.1029/95WR02966.

705 Leurgans, S. E., Moyeed, R. A. and Silverman, B. W., 1993: Canonical Correlation Analysis when the Data  
706 are Curves, *Journal of the Royal Statistical Society. Series B (Methodological)*, **55** (3), 725-740.

707 Lima, C. H. R. and Lall, U., 2010: Spatial scaling in a changing climate: A hierarchical bayesian model for  
708 non-stationary multi-site annual maximum and monthly streamflow, *Journal of Hydrology*, **383**  
709 (3–4), 307-318, doi: 10.1016/j.jhydrol.2009.12.045.

710 Lockwood, J. G., 1999: Is Potential Evapotranspiration and Its Relationship with Actual  
711 Evapotranspiration Sensitive to Elevated Atmospheric CO2 Levels?, *Climatic Change*, **41** (2), 193-  
712 212, doi: 10.1023/A:1005469416067.

713 MacDonald, G. M. and Case, R. A., 2005: Variations in the Pacific Decadal Oscillation over the past  
714 millennium, *Geophys. Res. Lett.*, **32** (8), L08703, doi: 10.1029/2005gl022478.

715 Matalas, N. C., Landwehr, J. M. and Wolman, M. G., 1982: Prediction in Water Management. *Scientific*  
716 *basis of water management*, National Academy, pp. 118-122.

717 McGowan, H. A., Marx, S. K., Denholm, J., Soderholm, J. and Kamber, B. S., 2009: Reconstructing annual  
718 inflows to the headwater catchments of the Murray River, Australia, using the Pacific Decadal  
719 Oscillation, *Geophys. Res. Lett.*, **36** (6), L06707, doi: 10.1029/2008gl037049.

720 Meko, D., Stockton, C. W. and Boggess, W. R., 1995: The tree-ring record of sever sustained drought,  
721 *JAWRA Journal of the American Water Resources Association*, **31** (5), 789-801, doi:  
722 10.1111/j.1752-1688.1995.tb03401.x.

723 Monteith, J., 1965: Evaporation and environment, *Symp. Soc. Exp. Biol*, **4**.

724 Nohara, D., Kitoh, A., Hosaka, M. and Oki, T., 2006: Impact of Climate Change on River Discharge  
725 Projected by Multimodel Ensemble, *Journal of Hydrometeorology*, **7** (5), 1076-1089, doi:  
726 10.1175/JHM531.1.

727 Nowak, K., Prairie, J., Rajagopalan, B. and Lall, U., 2010: A nonparametric stochastic approach for  
728 multisite disaggregation of annual to daily streamflow, *Water Resources Research*, **46** (8),  
729 W08529, doi: 10.1029/2009WR008530.

730 Palmer, J. G., Cook, E. R., Turney, C. S. M., Allen, K., Fenwick, P., Cook, B. I., O'Donnell, A., Lough, J.,  
731 Grierson, P. and Baker, P., 2015: Drought variability in the eastern Australia and New Zealand  
732 summer drought atlas (ANZDA, CE 1500–2012) modulated by the Interdecadal Pacific  
733 Oscillation, *Environmental Research Letters*, **10** (12), 124002, doi: 10.1088/1748-  
734 9326/10/12/124002.

735 Pederson, G. T., 2013, Multi-century perspectives on current and future streamflow in the Missouri River  
736 Basin, viewed 2015/09/08, <

737 Pederson, N., Bell, A. R., Cook, E. R., Lall, U., Devineni, N., Seager, R., Eggleston, K. and Vranes, K. P.,  
738 2012: Is an Epic Pluvial Masking the Water Insecurity of the Greater New York City Region?,  
739 *Journal of Climate*, **26** (4), 1339-1354, doi: 10.1175/jcli-d-11-00723.1.

740 Piechota, T. C. and Dracup, J. A., 1996: Drought and Regional Hydrologic Variation in the United States:  
741 Associations with the El Niño-Southern Oscillation, *Water Resources Research*, **32** (5), 1359-  
742 1373, doi: 10.1029/96WR00353.

743 Pizarro, G. and Lall, U., 2002: El Niño-induced Flooding in the U.S. West: What Can We Expect, *Eos Trans.*  
744 *AGU*, **82** (32), 349-352, doi: 10.1029/2002EO000255.

745 Prairie, J., Nowak, K., Rajagopalan, B., Lall, U. and Fulp, T., 2008: A stochastic nonparametric approach  
746 for streamflow generation combining observational and paleoreconstructed data, *Water*  
747 *Resources Research*, **44** (6), W06423, doi: 10.1029/2007wr006684.

748 Quinn, W. H., 1992: A study of the Southern Oscillation-related climatic activity for A.D 622-1990  
749 incorporating Nile River flood data. *El Niño: Historical and Paleoclimatic Aspects of the Southern*  
750 *Ocean*, Henry F. Diaz and Vera Markgraf, Eds., Cambridge University Press, pp. 119-149.

751 Rajagopalan, B. and Lall, U., 1999: A k-nearest-neighbor simulator for daily precipitation and other  
752 weather variables, *Water Resources Research*, **35** (10), 3089-3101, doi:  
753 10.1029/1999WR900028.

754 Redmond, K. T. and Koch, R. W., 1991: Surface Climate and Streamflow Variability in the Western United  
755 States and Their Relationship to Large-Scale Circulation Indices, *Water Resources Research*, **27**  
756 (9), 2381-2399, doi: 10.1029/91WR00690.

757 Routson, C. C., Woodhouse, C. A. and Overpeck, J. T., 2011: Second century megadrought in the Rio  
758 Grande headwaters, Colorado: How unusual was medieval drought?, *Geophys. Res. Lett.*, **38**  
759 (22), L22703, doi: 10.1029/2011gl050015.

760 Seager, R., Ting, M., Held, I., Kushnir, Y., Lu, J., Vecchi, G., Huang, H.-P., Harnik, N., Leetmaa, A., Lau, N.-  
761 C., Li, C., Velez, J. and Naik, N., 2007: Model Projections of an Imminent Transition to a More  
762 Arid Climate in Southwestern North America, *Science*, **316** (5828), 1181-1184, doi:  
763 10.1126/science.1139601.

764 Segond, M. L., Onof, C. and Wheater, H. S., 2006: Spatial-temporal disaggregation of daily rainfall from a  
765 generalized linear model, *Journal of Hydrology*, **331** (3-4), 674-689, doi:  
766 10.1016/j.jhydrol.2006.06.019.

767 Sheffield, J., Wood, E. F. and Roderick, M. L., 2012: Little change in global drought over the past 60  
768 years, *Nature*, **491** (7424), 435-438, doi: 10.1038/nature11575.

769 Smerdon, J. E., Cook, B. I., Cook, E. R. and Seager, R., 2015: Bridging Past and Future Climate across  
770 Paleoclimatic Reconstructions, Observations, and Models: A Hydroclimate Case Study, *Journal of*  
771 *Climate*, **28** (8), 3212-3231, doi: 10.1175/JCLI-D-14-00417.1.

772 Smith, S. R., Green, P. M., Leonardi, A. P. and O'Brien, J. J., 1998: Role of Multiple-Level Tropospheric  
773 Circulations in Forcing ENSO Winter Precipitation Anomalies, *Monthly Weather Review*, **126**  
774 (12), 3102-3116, doi: 10.1175/1520-0493(1998)126<3102:ROMLTC>2.0.CO;2.

775 Speer, J. H., 2010: *Fundamentals of tree-ring research*. University of Arizona Press, 333 pp.

776 St. George, S., 2010: Tree Rings as Paleoflood and Paleostage Indicators. *Tree Rings and Natural*  
777 *Hazards: A State-of-Art*, Markus Stoffel, Michelle Bollschweiler, R. David Butler, and H. Brian  
778 Luckman, Eds., Springer Netherlands, pp. 233-239.

779 Steinschneider, S. and Lall, U., 2015: Daily Precipitation and Tropical Moisture Exports across the Eastern  
780 United States: An Application of Archetypal Analysis to Identify Spatiotemporal Structure,  
781 *Journal of Climate*, **28** (21), 8585-8602, doi: 10.1175/JCLI-D-15-0340.1.

782 Stone, E. and Cutler, A., 1996: Introduction to archetypal analysis of spatio-temporal dynamics, *Physica*  
783 *D: Nonlinear Phenomena*, **96** (1-4), 110-131, doi: 10.1016/0167-2789(96)00016-4.

784 Thomas, D. M. and Benson, M. A., 1970: Generalization of streamflow characteristics from drainage-  
785 basin characteristics, Water Supply Paper, 1975, pp.

786 Thornthwaite, C. W., 1948: An Approach toward a Rational Classification of Climate, *Geographical*  
787 *Review*, **38** (1), 55-94, doi: 10.2307/210739.

788 Tierney, J. E., Oppo, D. W., Rosenthal, Y., Russell, J. M. and Linsley, B. K., 2010: Coordinated hydrological  
789 regimes in the Indo - Pacific region during the past two millennia, *Paleoceanography*, **25**,  
790 PA1102, doi: 10.1029/2009PA001871.

791 Tootle, G. A. and Piechota, T. C., 2006: Relationships between Pacific and Atlantic ocean sea surface  
792 temperatures and U.S. streamflow variability, *Water Resources Research*, **42** (7), n/a-n/a, doi:  
793 10.1029/2005WR004184.

794 Torrence, C. and Compo, G. P., 1998: A Practical Guide to Wavelet Analysis, *Bulletin of the American*  
795 *Meteorological Society*, **79** (1), 61-78, doi: 10.1175/1520-  
796 0477(1998)079<0061:APGTWA>2.0.CO;2.

797 Trenberth, K. E., Branstator, G. W. and Arkin, P. A., 1988: Origins of the 1988 North American Drought,  
798 *Science*, **242** (4886), 1640-1645, doi: 10.1126/science.242.4886.1640.

799 U.S. Geological Survey, 2011, GAGES-II: Geospatial Attributes of Gages for Evaluating Streamflow,  
800 viewed 06/19/2015,  
801 <[http://water.usgs.gov/GIS/metadata/usgswrd/XML/gagesII\\_Sept2011.xml](http://water.usgs.gov/GIS/metadata/usgswrd/XML/gagesII_Sept2011.xml)>.

802 US Army Corps of Engineers, 2006: Missouri River Mainstem Reservoir System - Master Water Control  
803 Manual - Missouri River Basin, pp. 431.

804 Valencia, D. and Schaake, J. C., 1973: Disaggregation processes in stochastic hydrology, *Water Resour.*  
805 *Res*, **9** (3), 580-585, doi: 10.1029/WR009i003p00580.

806 Vance, T. R., van Ommen, T. D., Curran, M. A. J., Plummer, C. T. and Moy, A. D., 2012: A Millennial Proxy  
807 Record of ENSO and Eastern Australian Rainfall from the Law Dome Ice Core, East Antarctica,  
808 *Journal of Climate*, **26** (3), 710-725, doi: 10.1175/jcli-d-12-00003.1.

809 Vinod, H. D., 1976: Canonical ridge and econometrics of joint production, *Journal of Econometrics*, **4** (2),  
810 147-166, doi: 10.1016/0304-4076(76)90010-5.

811 Vogel, R. M., Lall, U., Cai, X., Rajagopalan, B., Weiskel, P., Hooper, R. P. and Matalas, N. C., 2015:  
812 Hydrology: The interdisciplinary science of water, *Water Resources Research*, **51** (6), 4409-4430,  
813 doi: 10.1002/2015WR017049.

814 White, I. R., Royston, P. and Wood, A. M., 2011: Multiple imputation using chained equations: Issues and  
815 guidance for practice, *Statistics in Medicine*, **30** (4), 377-399, doi: 10.1002/sim.4067.

816 Wilson, R., Cook, E., D'Arrigo, R., Riedwyl, N., Evans, M. N., Tudhope, A. and Allan, R., 2010:  
817 Reconstructing ENSO: the influence of method, proxy data, climate forcing and teleconnections,  
818 *Journal of Quaternary Science*, **25** (1), 62-78, doi: 10.1002/jqs.1297.

819 Wise, E. K., 2010: Spatiotemporal variability of the precipitation dipole transition zone in the western  
820 United States, *Geophysical Research Letters*, **37** (7), doi: 10.1029/2009GL042193.

821 Woodhouse, C. A. and Meko, D., 1997: Number of Winter Precipitation Days Reconstructed from  
822 Southwestern Tree Rings, *Journal of Climate*, **10** (10), 2663-2669, doi: 10.1175/1520-  
823 0442(1997)010<2663:NOWPDR>2.0.CO;2.

824 Woodhouse, C. A. and Overpeck, J. T., 1998: 2000 years of drought variability in the central United  
825 States, *Bulletin of the American Meteorological Society*, **79** (12), 2693, doi.

826 Woodhouse, C. A., Gray, S. T. and Meko, D. M., 2006: Updated streamflow reconstructions for the Upper  
827 Colorado River Basin, *Water Resources Research*, **42** (5), W05415, doi: 10.1029/2005wr004455.

828 Woodhouse, C. A., Kunkel, K. E., Easterling, D. R. and Cook, E. R., 2005: The twentieth-century pluvial in  
829 the western United States, *Geophysical Research Letters*, **32** (7), doi: 10.1029/2005GL022413.

830 Woodhouse, C. A., Meko, D. M., MacDonald, G. M., Stahle, D. W. and Cook, E. R., 2010: A 1,200-year  
831 perspective of 21st century drought in southwestern North America, *Proceedings of the National*  
832 *Academy of Sciences*, **107** (50), 21283-21288, doi: 10.1073/pnas.0911197107.

833 Woodhouse, C. A., Lukas, J., Brice, B., Hirschboeck, K., Hartmann, H., Kostuk, M., Lay, E., Martinex, D.,  
834 McMahan, B. and Shah, A., 2002, Tree rings and streamflow, viewed 2015/06/01,  
835 <<http://treeflow.info/content/tree-rings-and-streamflow>>.

836

1 Can a paleo-drought record be used to  
2 reconstruct streamflow? A case-study for  
3 the Missouri River Basin

4 Michelle Ho<sup>1</sup>, Upmanu Lall<sup>1,2</sup>, Edward R. Cook<sup>3</sup>

5 1. Columbia Water Center, Columbia University, New York, NY, USA

6 2. Department of Earth and Environmental Engineering, Columbia University, New York, New York, USA

7 3. Lamont-Doherty Earth Observatory of Columbia University, Palisades, New York 10964, USA

8 Corresponding author: Michelle Ho. Email: mh3538@columbia.edu, Phone: 212-854-7081

9

10 Keywords: paleoclimate, regularized canonical correlation analysis, streamflow reconstruction

11 Key points

- 12 • Streamflow is reconstructed from an existing paleoclimate drought index
- 13 • Methodological innovation using rCCA addresses very high dimensional dataset
- 14 • Reconstructed streamflow provides insights into extreme, persistent high and low flow

## 15 Abstract

16 Recent advances in paleoclimatology have revealed dramatic long-term hydro-climatic variations that  
17 provide a context for limited historical records. A notable dataset derived from a relatively dense  
18 network of paleoclimate proxy records in North America is the Living Blended Drought Atlas (LBDA): a  
19 gridded tree-ring based reconstruction of summer Palmer Drought Severity Index. This index has been  
20 used to assess North American drought frequency, persistence and spatial extent over the past two  
21 millennia. Here, we explore whether the LBDA can be used to reconstruct annual streamflow. Relative to  
22 streamflow reconstructions that use tree rings within the river basin of interest, the use of a gridded  
23 proxy poses a novel challenge. The gridded series have high spatial correlation, since they rely on tree  
24 rings over a common radius of influence. A novel algorithm for reconstructing streamflow using  
25 regularized canonical regression and inputs of local and global covariates is developed and applied over  
26 the Missouri River Basin, as a test case. Effectiveness in reconstruction is demonstrated with  
27 reconstructions showing periods where streamflow deficits may have been more severe than during  
28 recent droughts (e.g. the Civil War, Dust Bowl and 1950s droughts). The maximum persistence of  
29 droughts and floods over the past 500 years far exceed those observed in the instrumental record and  
30 periods of multi-decadal variability in the 1500s and 1600s are detected. Challenges for an extension to  
31 a national streamflow reconstruction or applications using other gridded paleoclimate datasets such as  
32 adequate spatial coverage of streamflow and applicability of annual reconstructions are discussed.

Michelle Ho 4/6/2016 8:25 AM

Deleted: A

Michelle Ho 4/6/2016 8:25 AM

Deleted: in the past few decades

Michelle Ho 4/6/2016 8:14 AM

Deleted:

Michelle Ho 4/6/2016 8:25 AM

Deleted: the

Michelle Ho 4/6/2016 8:19 AM

Deleted: annually-resolved

Michelle Ho 4/6/2016 8:15 AM

Deleted: the

Michelle Ho 4/6/2016 8:25 AM

Deleted: (June-August)

Michelle Ho 4/6/2016 8:21 AM

Deleted: (PDSI)

Michelle Ho 4/6/2016 8:15 AM

Deleted: the

Michelle Ho 4/6/2016 8:15 AM

Deleted: of North American droughts

Michelle Ho 4/6/2016 8:27 AM

Deleted: series to reconstruct watershed processes such as streamflow

Michelle Ho 4/6/2016 8:18 AM

Deleted: o

Michelle Ho 4/6/2016 8:19 AM

Deleted: for streamflow reconstruction

Michelle Ho 4/6/2016 8:28 AM

Deleted: streamflow

Michelle Ho 4/6/2016 8:17 AM

Deleted: in recent history

Michelle Ho 4/6/2016 8:13 AM

Deleted: An analysis of drought and flood epochs using a

Michelle Ho 4/6/2016 8:13 AM

Deleted: simple threshold analysis shows that t

Michelle Ho 4/6/2016 8:14 AM

Deleted: . A wavelet analysis reveals

Michelle Ho 4/6/2016 8:24 AM

Deleted: streamflow

## 55 1 Introduction

56 The concern with anthropogenic climate change and its hydrologic impacts has focused interest on how  
57 long term climate variability may impact streamflow (e.g. Nohara et al., 2006; Seager et al., 2007). The  
58 consequence of using short records to “over-allocate” the flows of major rivers are often cited as an  
59 example of the need for long records that can better inform the possible range of long-term variations  
60 of streamflow (Tootle and Piechota, 2006; McGowan et al., 2009; Woodhouse et al., 2010). Continuous  
61 records of streamflow in the US span several decades at best. Advances in paleoclimatology in the past  
62 few decades have provided opportunities across the world to extend the range of hydroclimatic  
63 variability (e.g. Quinn, 1992; Jones and Mann, 2004; Tierney et al., 2010; Gallant and Gergis, 2011; Vance  
64 et al., 2012; Cook et al., 2013b; Devineni et al., 2013; Ho et al., 2015). While considerable uncertainty  
65 clouds the projections of hydroclimatic states towards the end of the 21<sup>st</sup> century, in the near to  
66 medium term paleoclimate information may be crucial to inform the interannual to decadal variability of  
67 regional water availability as indicated by streamflow for reservoir operation, and agricultural and other  
68 water use decisions.

69 Paleoclimate reconstructions have been developed using proxies that typically span the past 1000-2000  
70 years (also known as the Common Era). The North American region has a relatively dense network of  
71 high-resolution paleoclimate proxy records, primarily comprised of tree-ring chronologies. Tree-ring-  
72 proxy records have been used to assess various components of environmental variations (Fritts, 1976)  
73 including drought severity (e.g. Routson et al., 2011; Cook et al., 2015a), pluvials (e.g. Woodhouse et al.,  
74 2005; Pederson et al., 2012), streamflow variability (e.g. Woodhouse et al., 2006; Prairie et al., 2008;  
75 Allen et al., 2013; Devineni et al., 2013), and precipitation frequency (Woodhouse and Meko, 1997) in  
76 addition to enabling comparisons of past climate with projected climate scenarios (e.g. Ault et al., 2014;  
77 Cook et al., 2015a; Smerdon et al., 2015).

Michelle Ho 3/6/2016 6:01 PM

Deleted: the Colorado River is

79 Studies focused on the reconstruction of tree-ring-based paleo-hydrology (e.g. annual and season  
80 streamflow and floods) typically utilize proxies derived from tree-ring networks within or near the  
81 catchment region as predictors (e.g. Woodhouse et al., 2006; St. George, 2010; Pederson et al., 2012;  
82 Devineni et al., 2013). However, these networks are spatially irregular with record lengths varying across  
83 chronologies. An alternative to using spatially and temporally irregular tree-ring chronologies as model  
84 predictors is to use an existing derivative of these records, namely the Living Blended Drought Atlas  
85 (LBDA) (Cook et al., 2010a). The LBDA is a paleoclimate reconstruction of the summer (June-August)  
86 Palmer Drought Severity Index (PDSI) that is gridded across North America on a  $0.5^\circ \times 0.5^\circ$   
87 latitude/longitude grid with reconstructions dating back as far as 2000 years. These records are  
88 temporally complete over the Conterminous United States (CONUS) from 1473 onward. The LBDA, or its  
89 predecessor the North American Drought Atlas (Cook and Krusic, 2004), have been used to assess the  
90 frequency and spatial distribution of droughts over the past millennia (e.g. Herweijer et al., 2007; Cook  
91 et al., 2013a).

## 92 2 A proposal for using PDSI to reconstruct streamflow

93 The intent of the modeling case study presented here is to develop a suitable framework with which  
94 streamflow within the CONUS may be reconstructed using a tree-ring-based reconstruction of the PDSI.  
95 In developing the modelling framework we consider; 1) the constraints posed by the LBDA and  
96 implications for reconstructing streamflow; 2) possible temporal resolutions (monthly, seasonal or  
97 annual) for direct streamflow reconstruction using the LBDA data; 3) how to best use local and far-field  
98 LBDA information for local streamflow reconstruction; and 4) provides insights from the 500 year  
99 reconstruction of the multi-site Missouri River Basin flows as to the decadal and longer variability of  
100 streamflow in the region. Given the existence of the spatially and temporally complete LBDA record over  
101 the last 500 years covering the CONUS, we explore whether the LBDA, a reconstruction of PDSI using

Michelle Ho 3/6/2016 6:13 PM

Deleted: s

103 | tree-ring chronologies, could be a reasonable predictor of streamflow variability. In this case the variable  
104 | to be reconstructed is annual streamflow, with the aim of eventually reconstructing paleoclimate records  
105 | of streamflow across the CONUS, an undertaking that has not previously been attempted using the  
106 | LBDA or tree-rings.

107 | The motivation for using the LBDA stems from our understanding that the growth of moisture-limited  
108 | trees, from which the LBDA is derived, are in part governed by climatic forcings that drive soil moisture  
109 | availability. That is, given a vector of unspecified climate variables,  $\mathbf{C}_t$  (where  $\mathbf{C}$  may be comprised of, but  
110 | not limited to, climate variables such as temperature, rainfall, wind, soil moisture and radiation) we can  
111 | define PDSI as  $PDSI_t = f_1(\mathbf{C}_t)$ . Streamflow is also a derivative of a number of climate variables and can  
112 | similarly be defined as  $Q_t = f_2(\mathbf{C}_t)$ . We seek to determine if it is possible to derive and fit a function  $f_3$  that  
113 | relates streamflow to PDSI where  $Q_t = f_3(PDSI_t)$ . Where suitable instrumental records of climate are  
114 | available, we may derive  $f_1$  and this has been approximated in part using methods such as the  
115 | Thornthwaite potential evapotranspiration (PET) equation (Thornthwaite, 1948) and the Penman-  
116 | Monteith PET equation (Monteith, 1965) to varying degrees of success (Lockwood, 1999; Sheffield et al.,  
117 | 2012; Dai, 2013). Consideration of the joint probability distributions  $f(\mathbf{C}_t | PDSI_t)$  and  $f(Q_t, PDSI_t)$  enables  
118 | the vector of climate variables to be reconstructed by capitalizing on the climate-PDSI relationship and  
119 | the tree-ring chronology-PDSI relationship given  $f(\mathbf{C}_t | PDSI_t) f(PDSI_t | \text{tree-ring chronology}_t)$ . Paleoclimate  
120 | streamflow could similarly be derived by implementing  $f_2$ , namely  $f(Q_t | \mathbf{C}_t) f(\mathbf{C}_t | PDSI_t) f(PDSI_t | \text{tree-ring}$   
121 | ***chronology}\_t***). However, the challenge in this approach is the selection of an appropriate vector of  
122 | climate variables, many of which may be sparsely observed or unobserved in the instrumental record.  
123 | Therefore, we consider modelling paleoclimate streamflow through  $f(Q_t | PDSI_t) f(PDSI_t | \text{tree-ring}$   
124 | ***chronology}\_t***). Since the reconstructed PDSI<sub>t</sub> in the LBDA is really  $E[PDSI_t | \text{tree-ring chronology}_t]$ , we start  
125 | by considering  $f(Q_t | PDSI_t)$ , where  $Q_t$  is the streamflow at one or more locations in a river basin, and  
126 |  $PDSI_t$  represents a vector of LBDA values at the gridded locations of the LBDA that can be a potential

Michelle Ho 3/6/2016 5:23 PM

Deleted: ,

Michelle Ho 3/6/2016 5:23 PM

Deleted: i

Michelle Ho 3/6/2016 5:23 PM

Deleted: ,

130 predictor of the streamflows. The relationship can be developed using contemporaneous values of  $PDSI_t$   
131 and historical  $Q_t$ . Subsequently, we can apply this relationship to the paleo estimates of  $PDSI_t$   
132 recognizing that we could use the expected values of  $PDSI_t$  reported in the LBDA, or simulations from  
133 the uncertainty distributions of  $PDSI_t$  reported in the LBDA, as conditioning variables to derive  
134 sequences of  $Q_t$ .

135 A key motivation for using the LBDA is that it is spatially complete across CONUS over the past 500 years  
136 and a successful streamflow reconstruction would have significant value in assessing national water  
137 planning and use strategies and in investigating the different temporal and spatial structures in  
138 streamflow, which differ from the LBDA (see Figures S1 and S2 in supporting information). Relative to  
139 reconstructions that use tree rings within the river basin of interest, the use of a gridded proxy series to  
140 reconstruct watershed processes such as streamflow poses a novel challenge. The gridded LBDA series  
141 have high spatial correlation, since they rely on tree rings over a common radius of influence around  
142 each grid point, nominally 450 km in this case or roughly the correlation decay  $e$ -folding distance  
143 between grid points of the instrumental PDSI. A novel algorithm for reconstruction that uses local and  
144 global covariates for streamflow reconstruction using regularized canonical regression is developed and  
145 applied over the Missouri River Basin, as a test case. This provides a proof of concept at a sub-  
146 continental scale as to whether the approach is feasible. The Missouri River Basin was chosen as it  
147 contains the only major river headwaters in the western US where extensive reconstructions of  
148 paleoclimate hydrology have not been undertaken (Driscoll, 2013) and also parallels current efforts to  
149 further develop tree-ring chronologies and streamflow reconstruction models using tree rings in the  
150 region (Pederson, 2013). A description of the case study region and data are presented in the following  
151 section while Section 4 presents initial diagnostics that inform the modeling approach developed in  
152 Section 5. Section 5 details how the above proposal of using PDSI to reconstruct streamflow is  
153 implemented in the Missouri River Basin. Section 6 provides model verification results and summaries of

154 the key modes of variability in the 500 year mean streamflow reconstructions of Missouri River Basin  
155 streamflow. The final section reviews outstanding questions as to modeling uncertainties, and the  
156 challenges for extending the model presented here to a national scale.

## 157 3 Case Study Region and Data

### 158 3.1 The Missouri River Basin

159 The Missouri River is the longest river and the second largest river basin in the US, draining an area over  
160 1.3 million km<sup>2</sup> that spans the southern portions of two Canadian provinces and ten states in the US  
161 (Missouri River basin boundary and key tributaries shown in Figure 1). The headwaters, which are largely  
162 snow-melt fed, are located in the Northern Rocky Mountains. These waters then flow through a largely  
163 semi-arid region to its confluence with the Mississippi River near St. Louis, Missouri (Galat et al., 2005).  
164 Land use in the Missouri River Basin is dominated by agricultural activities including cropping and  
165 grazing which cover 95% of the region (US Army Corps of Engineers, 2006), while the remaining land is  
166 used for recreation, transport, urban and industrial use including mining and energy sector activities  
167 (Galat et al., 2005). An improved perspective of streamflow variability would therefore be of benefit to  
168 the region in terms of managing and balancing the demands for water amongst the various sectors and  
169 users.

170 Precipitation and moisture availability in the Missouri River Basin are characterized by high precipitation  
171 in the western mountainous region, which averages over 1000 mm/year to the drier region in the rain  
172 shadow east of the Rocky Mountains where average annual precipitation is less than 400 mm/year.  
173 Precipitation increases towards the far eastern regions of the Missouri River Basin (Kunkel et al., 2013).  
174 Winter precipitation in the northern Missouri River Basin and in the mountainous regions to the west is  
175 related to the El Niño Southern Oscillation (ENSO) signal from the preceding summer and autumn

176 (Redmond and Koch, 1991). El Niño teleconnections typically manifest as upper level anticyclonic high  
177 pressure cells over the northwestern US and result in the northward displacement or splitting of the jet  
178 stream and anomalously dry conditions in the Missouri River Basin (Trenberth et al., 1988; Dettinger et  
179 al., 1998; Smith et al., 1998). Conversely, La Niña events typically result in wetter conditions in this  
180 region. ENSO impacts are modulated by the decadal scale variability in the northern Pacific Ocean with  
181 warm decadal phases resulting in a deep Aleutian low and corresponding ridging over the western US  
182 thereby enhancing El Niño conditions (Gershunov and Barnett, 1998; Brown and Comrie, 2004). The  
183 enhancement of La Niña impacts by a cool phase in the northern Pacific decadal signature is particularly  
184 noticeable in the northern Missouri River Basin (Wise, 2010). Both Pacific and Atlantic Ocean influences  
185 are seen in the Northern Great Plains with the Great Plains low level jet, originating from the Gulf of  
186 Mexico, enhancing summer precipitation in this region (Higgins et al., 1997).

### 187 3.2 Streamflow data

188 Monthly streamflow data for the Missouri River Basin were obtained from the United States Geological  
189 Survey (USGS) Surface-Water Daily Data for the Nation (<http://nwis.waterdata.usgs.gov/nwis/sw>). The  
190 streamflow data are from stations included in the USGS's GAGES-II network (U.S. Geological Survey,  
191 2011) within the Missouri River Basin boundary and are reference gauges identified by the USGS as the  
192 least-disturbed watersheds with minimal regulation. The selected gauges meet a criterion of data  
193 spanning 40 years with less than 5% missing data (Figure 1, Table S1) and results in 55 streamflow  
194 gauges, 46 of which are also in the USGS Hydro Climatic Data Network.

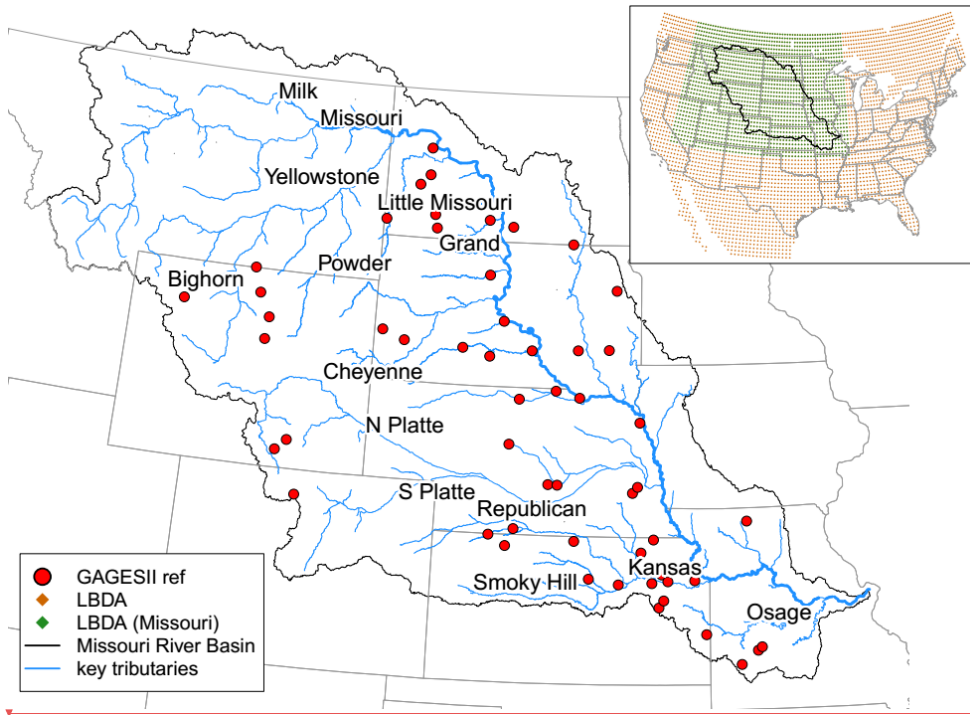


Figure 1. The Missouri River Basin, key tributaries, locations of gauged streamflow used in the analysis, and (inset) the LBDA grids used in the analysis. LBDA grids in the Missouri River Basin region shown in green.



195

196  
197

198 Three of the selected stations had missing monthly values. These values were imputed using multiple  
 199 imputation by chained equations (MICE) and a method of predictive mean matching (Buuren and  
 200 Groothuis-Oudshoorn, 2011). Monthly streamflow imputation was conducted using streamflow from  
 201 the three closest stations and a cosine function to represent a seasonal signal. The number of  
 202 repetitions in MICE was determined using a rule of thumb method proposed by White et al. (2011) and  
 203 the multiple imputed monthly values were averaged across the repetitions.

204 Monthly data was aggregated into annual streamflow data using a calendar year instead of a water year  
 205 because the average driest month of streamflow at most stations occurred in either December or  
 206 January. Start and end years with incomplete data were excluded. The resulting annual data spans from

208 1929 to 2014 with annual record lengths varying between 39 and 85 years after aggregation. Streamflow  
209 was logarithmically transformed since that leads to a nearly Gaussian distribution for annual streamflow.  
210 Annual streamflow records of zero were replaced with half the minimum annual streamflow prior to the  
211 application of the log transform.

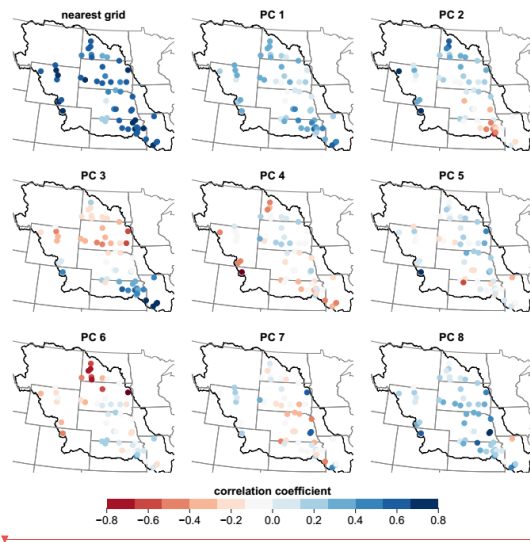
### 212 3.3 A Paleoclimate Record of North American Drought

213 The LBDA is an updated version of the seminal North American Drought Atlas (NADA) (Cook et al., 1999;  
214 Cook and Krusic, 2004; Cook et al., 2010a), which is a paleoclimate reconstruction of the summer (June  
215 to August – JJA) PDSI based on a network of tree-ring chronologies. The LBDA has a spatial resolution of  
216  $0.5^\circ \times 0.5^\circ$  latitude/longitude and incorporates information from 1845 tree-ring chronologies, an  
217 improvement over previous NADA versions that were informed by fewer tree-ring chronologies and  
218 were calculated over coarser grids (Cook et al., 1999; Cook et al., 2004). The LBDA is spatially complete  
219 over the CONUS region from 1473-2005 and includes instrumental data from 1979 onwards. A region of  
220 the LBDA ranging from  $23^\circ\text{N}$  to  $52^\circ\text{N}$  and  $125^\circ\text{W}$  to  $66^\circ\text{W}$  (see green dots in Figure 1) was extracted for  
221 the analysis. This broad region extending beyond the CONUS region was selected to capture patterns of  
222 LBDA variability relevant to the CONUS. A comparative analysis between an instrumental-based gridded  
223 PDSI dataset and the LBDA showed that they are highly correlated, with correlation values significant at  
224 the 99% level and similar variance (results not shown here). The instrumental LBDA data post-1978 were  
225 therefore also included in the modeling analysis performed here.

## 226 4 Initial Diagnostic Analyses

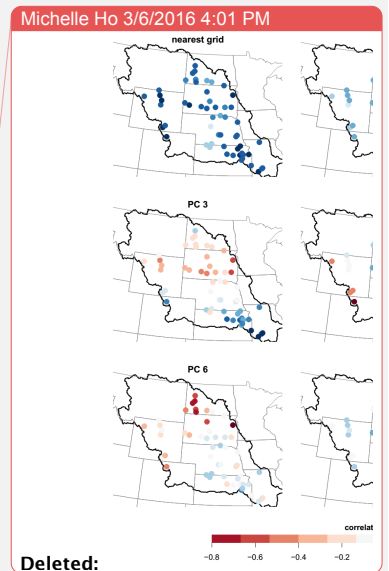
227 Initial diagnostic analyses were performed using both parametric and non-parametric correlations  
228 (Pearson and Kendall correlations respectively) between the LBDA and the log-normal streamflow  
229 series. The two correlations measures yielded similar results suggesting a near-Gaussian linear  
230 dependency (Pizarro and Lall, 2002). Different levels of temporal aggregation were tested including

231 monthly, rolling seasonal, bi-seasonal and annual streamflow. Different representations of the LBDA  
 232 were also tested including using the LBDA grid located nearest to each streamflow gauge, using LBDA  
 233 grids surrounding each streamflow gauge within a given diameter, and principal components (PCs)  
 234 (Jolliffe, 2002) and archetype analysis (Cutler and Breiman, 1994) of US-wide and Missouri River Basin  
 235 region LBDA (see orange and green diamonds respectively in Figure 1). An annual temporal aggregation  
 236 of streamflow was found to produce the strongest signal (Figure 2 showing correlation results between  
 237 streamflow gauges in the Missouri River Basin and the nearest LBDA grid and PCs of US-wide LBDA).  
 238 Diagnostic tests using the nearest LBDA grid were superior and reflects the ability of the point-by-point  
 239 regression method used for the LBDA to preserve local climate details in the PDSI reconstructions that  
 240 are also related to streamflow.



241  
 242 *Figure 2. Pearson correlation between annual streamflow and a) closest LBDA grid b)-i) PC1 – PC8 of US-wide LBDA*

243 High correlations were also noted for some streamflow stations with LBDA in the surrounding region  
 244 (see supporting information) particularly for streamflow records with weaker correlations with the  
 245 nearest grid (e.g. gauges in south-central North Dakota and on the border of Nebraska and Kansas). We



247 therefore considered information from LBDA grids within a 450 km radius consistent with the tree  
248 chronology search radius used to form the LBDA. Furthermore, given that large-scale climate and  
249 weather patterns influence local hydroclimatic conditions (Woodhouse et al., 2002), one needs to also  
250 consider the relationships with these larger modes of variability.

251 Correlations between streamflow and the first eight PCs of LBDA grids across the US (LBDA locations  
252 shown in Figure 1, correlations in Figure 2 and PC loading patterns shown in Figure 3) show that the  
253 Missouri River Basin streamflow is correlated with PC1, a PC representing overall US LBDA variability.  
254 However, these correlations are weak in comparison with correlations with the nearest LBDA grid. The  
255 north/south loading pattern of PC 2 is reflected in the change in correlation sign for stations north and  
256 south of the Nebraska and Kansas border. Similarly, the east-west difference in the loadings of PC3 leads  
257 to strong positive correlations in the downstream reaches of the Basin in Kansas and Missouri and  
258 negative correlations in Nebraska and Wyoming. The correlation results between streamflow and PCs of  
259 CONUS LBDA suggest that the large-scale modes of variability could provide additional information for  
260 modelling streamflow. Further details of the eight PCs of the CONUS LBDA and selection methods are  
261 provided in the supporting information.

## 262 5 Modeling approach and performance metrics

263 Here, we present a suite of plausible methods that could be implemented to model streamflow using  
264 the LBDA as the covariate(s) (Section 5.1). We justify the selection of a model that incorporates the use  
265 of regularized canonical correlation (rCCA), which is further described in Section 5.2. A description of the  
266 model performance metrics is provided in Section 5.3.

267 5.1 Model design and preliminary model assessment

268 We seek to use log-linear models to quantify the relationship between streamflow at individual gauges  
269 and a suite of site-specific LBDA records to facilitate a paleoclimate reconstruction of streamflow in the  
270 Missouri River Basin. Keeping in mind the application to streamflow record extension using the  
271 reconstructed PDSI values available in the LBDA we consider the following steps:

272 1. *Instrumental Period (32-77 years at 55 locations)*: Use LBDA reconstructed PDSI to estimate the  
273 relationship between log transformed annual streamflow and a selection of LBDA records specific to the  
274 target streamflow site, namely  $f(\ln(Q_t) | PDSI_t)$ , where  $PDSI_t$  is really the expected value of PDSI informed  
275 by the tree-ring chronologies,  $E[PDSI_t | \text{tree-ring chronology}_t]$ , from 1929 to 2005 to maximize the  
276 overlap between the two sets of variables.

277 2. *Paleo Period (1473 to 1929)*: Use  $f(\ln(Q_t) | PDSI_t)$  estimated in the previous step with the LBDA  
278 reconstructed PDSI,  $E[PDSI_t | \text{tree-ring chronology}_t]$ , to estimate  $\ln(Q_t)$  prior to 1929 (estimates of  $\ln(Q_t)$   
279 post 1929 will also be shown for comparison).

280 Several reconstruction model designs were considered given the initial diagnostic results. These  
281 included:

- 282 1) Developing one model for each streamflow station with predictors comprised of either:
- 283 a. the LBDA grid located closest to the streamflow gauge;
  - 284 b. PCs of the LBDA from a region around the Missouri River Basin (Figure 1, green diamonds)  
285 within a multilinear model framework;
  - 286 c. canonical correlation analysis with regularization (rCCA described in Section 5.2) using LBDA  
287 grids within a 450km radius; or
  - 288 d. rCCA using LBDA grids within a 450km radius in addition to the retained PCs of CONUS-wide  
289 LBDA.

- 290 2) Developing one model for all streamflow stations with predictors comprised of either:
- 291 a. PCs of both streamflow and CONUS-wide LBDA; or
- 292 b. rCCA using streamflow and LBDA information from either the Missouri River Basin region or the
- 293 CONUS region.

294 The fine resolution, gridded LBDA data results in a large number of highly correlated, potential

295 predictors in all of the model designs considered with the exception of 1a. Given the  $0.5^\circ \times 0.5^\circ$

296 resolution of the LBDA data, if a reconstruction of streamflow at a single gauge is considered using only

297 a 450 km radius of surrounding LBDA data, one has around 300 potential predictors including the

298 leading eight CONUS LBDA PCs. The length of annual streamflow records in the Missouri River Basin

299 ranges from 39 to 85 years, and thus it is clear that the number of potential predictors is significantly

300 greater than the number of observations available to fit the model. Of course, if one were to consider a

301 model with a simultaneous reconstruction of all the streamflow records, the dimension of the

302 estimation problem becomes even more challenging. Consequently, this is the first issue considered in

303 model development.

304 A single reconstruction model would be advantageous in its ability to consider all streamflow stations at

305 once rather than fitting 55 individual models. However, the varying streamflow record lengths and

306 differences in record periods would have resulted in large uncertainties as a large degree of annual

307 streamflow imputation would have been required. As a result, model designs under option 2 were not

308 further examined.

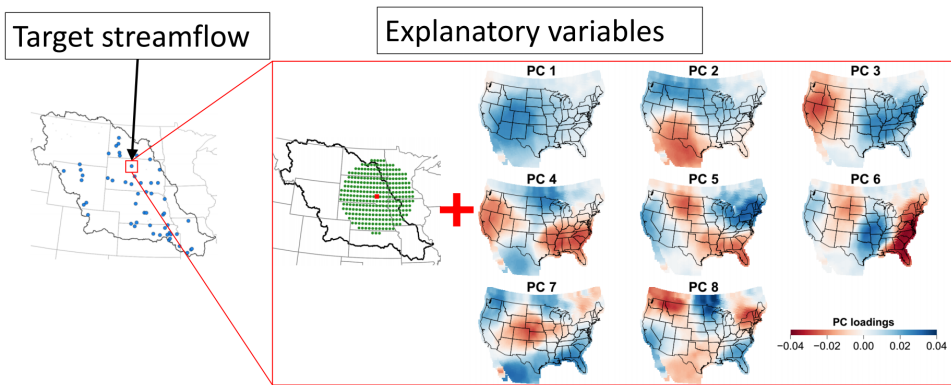
309 A cursory comparison of the viable models was conducted by comparing the coefficient of

310 determination for each fitted model. The development of individual models for each streamflow gauge

311 using either model 1a (the closest LBDA grid) or 1b (PCs of the LBDA in a region around the Missouri

312 River Basin) resulted in acceptable models of streamflow (mean adjusted  $R^2$  across the two models were

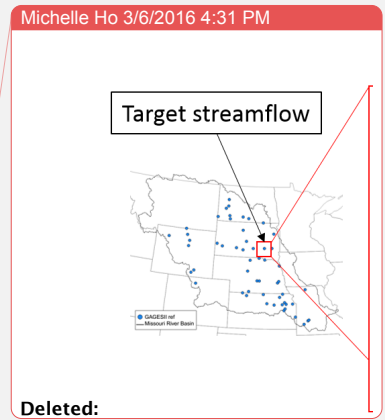
313 0.37 and 0.33 respectively). However, the spatially complete LBDA presents the opportunity to include a  
 314 wider variety of local and large-scale information and these models (model designs 1c and 1d) were  
 315 superior to the simpler models (model designs 1a and 1b). In addition, rCCA using both local and  
 316 CONUS-wide information (model 1d) resulted in slightly improved model results (mean adjusted  $R^2$  of  
 317 0.74 across all individual station models) over using only local information (model 1c, mean adjusted  $R^2$   
 318 of 0.71). We therefore selected model 1d (rCCA using LBDA grids within a 450km radius in addition to  
 319 the retained PCs of CONUS-wide LBDA) for further analysis. A schematic of the data used to fit this  
 320 model is shown in Figure 3, while a more detailed description of rCCA is provided in Section 5.2.



321  
 322 *Figure 3. Schematic of LBDA information to be included in the reconstruction model of one streamflow station in the Missouri*  
 323 *River Basin. The model inputs are the LBDA within a 450km radius and the first 8 PCs of US-wide LBDA.*

324 **5.2 Regularized canonical correlation analysis**

325 Methods such as principal component analysis (PCA, Jolliffe, 2002), archetype analysis (AA, Cutler and  
 326 Breiman, 1994; Stone and Cutler, 1996; Steinschneider and Lall, 2015) and canonical correlation analysis  
 327 (CCA, Hotelling, 1936) are typically used for dimension reduction in this setting. Given that the number  
 328 of potential predictors exceeds the number of observations, their high mutual correlation (>0.95 for  
 329 many of the adjacent grids), and an interest in exploring a multivariate streamflow response, we explore  
 330 the use of regularized canonical correlation where the regularization procedure is akin to ridge



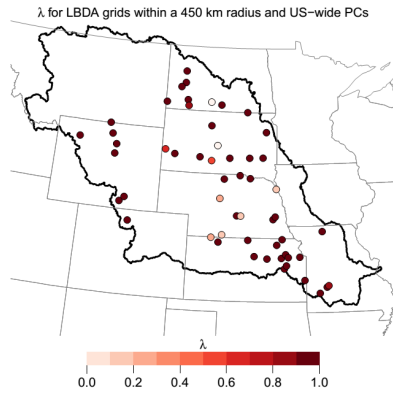
332 regression (De Bie and De Moor, 2003). Regularized canonical correlation was selected for use here to  
 333 capitalize on its ability to tailor the dimension reduction process to maximize the correlation between  
 334 the explanatory and target variables (in contrast to PCA where the explanatory and target variable  
 335 relationship is not considered).

336 Canonical correlation was introduced by Hotelling (1936) as a method of linearly transforming two  
 337 vector variables to canonical form to maximize the correlation between them. Consider two sets of  
 338 random variables represented by two matrices  $\mathbf{X}$  and  $\mathbf{Y}$ .  $\mathbf{X}$  is a  $n \times p$  matrix where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_p]$  containing  
 339  $n$  observations at  $p$  different locations, while  $\mathbf{Y}$  is a  $n \times q$  matrix where  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_q]$  containing  $n$   
 340 observations at  $q$  different locations. Both  $\mathbf{X}$  and  $\mathbf{Y}$  have finite variance matrices represented by  $\Sigma_{XX}$  and  
 341  $\Sigma_{YY}$  respectively. The covariance matrix between  $\mathbf{X}$  and  $\mathbf{Y}$  is  $\Sigma_{XY}$  while the covariance matrix between  $\mathbf{Y}$   
 342 and  $\mathbf{X}$  is  $\Sigma_{YX}$ . In this application,  $\mathbf{X}$  consists of the LBDA inputs specific to each streamflow gauge and  $\mathbf{Y}$  is  
 343 the streamflow at one station (i.e.  $p$  is large and  $q = 1$ ). Canonical correlation analysis involves rotating  
 344 the coordinate axes of both  $\mathbf{X}$  and  $\mathbf{Y}$  to new coordinate systems in order to clearly exhibit correlation  
 345 between  $\mathbf{X}$  and  $\mathbf{Y}$ . An arbitrary linear combination could be  $\mathbf{U} = \mathbf{X} \times \boldsymbol{\alpha}$  and  $\mathbf{V} = \mathbf{Y} \times \boldsymbol{\gamma}$  such that the correlation  
 346 between  $\mathbf{U}$  and  $\mathbf{V}$  is maximized.  $\mathbf{U}$  and  $\mathbf{V}$  yield the first pair of canonical variates, while  $\boldsymbol{\alpha}$  and  $\boldsymbol{\gamma}$  are the  
 347 vectors of canonical weights of length  $p$  and  $q$  respectively. This process may be repeated subject to the  
 348 constraint that following pairs of canonical variates are orthogonal to previous pairs with a maximum of  
 349  $\min(p, q)$  pairs obtained. In this application, the model of streamflow is applied station by station and  
 350 therefore only one pair of canonical variates are calculated and the first canonical variate of the LBDA is  
 351 used to fit the model of streamflow.

352 The correlation of successive pairs of canonical variates can be found using an eigen decomposition of  
 353  $\Sigma_{XX}^{-1} \Sigma_{XY} \Sigma_{YY}^{-1} \Sigma_{YX}$  and  $\Sigma_{YY}^{-1} \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}$ . The resulting  $\min(p, q)$  eigenvalues are common to both and the  
 354 square root of the eigenvalues yield the canonical correlation. The eigenvectors of  $\Sigma_{XX}^{-1} \Sigma_{XY} \Sigma_{YY}^{-1} \Sigma_{YX}$  and  
 355  $\Sigma_{YY}^{-1} \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}$  respectively yield the canonical weights  $\boldsymbol{\alpha}$  and  $\boldsymbol{\gamma}$  that are used to transform  $\mathbf{X}$  and  $\mathbf{Y}$ .

356 These weights are akin to the beta values in a multiple linear regression. Transforms using the  $i^{th}$   
357 eigenvector result in correlations corresponding to the square root of the  $i^{th}$  eigenvalue. Here, the  
358 weight for  $Y$  (streamflow) is 1 and the weights for  $X$  (LBDA) are given by the first eigenvector of  $\Sigma_{YY}^{-1}$   
359  $\Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}$ .

360 We employ regularization of the CCA process to address the issue of a large number of predictors  
361 relative to the number of observations (Vinod, 1976). Regularization is a smoothing process where a  
362 “roughness penalty”, also known as the regularization parameter ( $\lambda$ ), is introduced by converting  $\Sigma_{YY}$  and  
363  $\Sigma_{XX}$  to  $\Sigma_{YY} + \lambda_y I$  and  $\Sigma_{XX} + \lambda_x I$  respectively (Leurgans et al., 1993) and is similar to the technique of ridge  
364 regression (De Bie and De Moor, 2003). Values of  $\lambda$  range between 0 and 1 with larger  $\lambda$  values  
365 indicating a higher degree of smoothing. Regularization also enables the process of matrix inversion to  
366 be stabilized. A suitable value of  $\lambda$  is determined using a leave one out cross validation score, where  $\lambda_x$   
367 and  $\lambda_y$  are selected such that the correlation between the transformed datasets are maximized while the  
368 degrees of freedom used (as defined by Dijkstra, 2014) are limited to a maximum of  $n-10$ . If the criteria  
369 for the maximum degrees of freedom could not be achieved  $\lambda$  was set to one, otherwise  $\lambda$  was  
370 evaluated to two significant figures. No regularization was required for the single variable streamflow,  
371 whilst the LBDA was heavily regularized in almost all cases (Figure 4). rCCA was executed in R using the R  
372 package ‘CCA’ by González et al. (2008) and is freely available from the Comprehensive R Archive  
373 Network (CRAN <https://cran.r-project.org/>).



374

375 *Figure 4. CCA regularization parameter values for the explanatory variables (LBDA grids within a 450km radius and CONUS-wide*  
 376 *PCs) for each model of LBDA and streamflow.*

377 **5.3 Model performance metrics**

378 The performance of the model selected for further analysis includes verification using a leave-k-out  
 379 cross-validation procedure. Ten percent of the data (between 3 and 8 years out of a total of 32 and 77  
 380 years of streamflow data overlapping the LBDA record) are randomly selected and withheld from model  
 381 fitting. These values are then predicted from the model fit to the balance of the data. The entire process  
 382 is repeated 100 times, thus providing a set of 100 k-fold cross-validation samples. A comparison of the  
 383 model residuals resulting from both calibrated and verified model inputs is made for the 100 cross-  
 384 validation samples. The coefficient of efficiency (CE) and the reduction of error (RE) (Cook et al., 1994;  
 385 Wilson et al., 2010) are additional metrics calculated to verify the model. Both CE and RE are similar to  
 386 the Nash-Sutcliffe efficiency test, however, the metrics are normalized using the mean of the verification  
 387 and calibration data respectively. Namely, CE is defined as

388

$$CE = 1 - \frac{\sum_{i=1}^n (x_i - \hat{x}_i)}{\sum_{i=1}^n (x_i - \bar{x}_v)}$$

389 while RE is defined as

$$RE = 1 - \frac{\sum_{i=1}^n (x_i - \hat{x}_i)}{\sum_{i=1}^n (x_i - \bar{x}_c)}$$

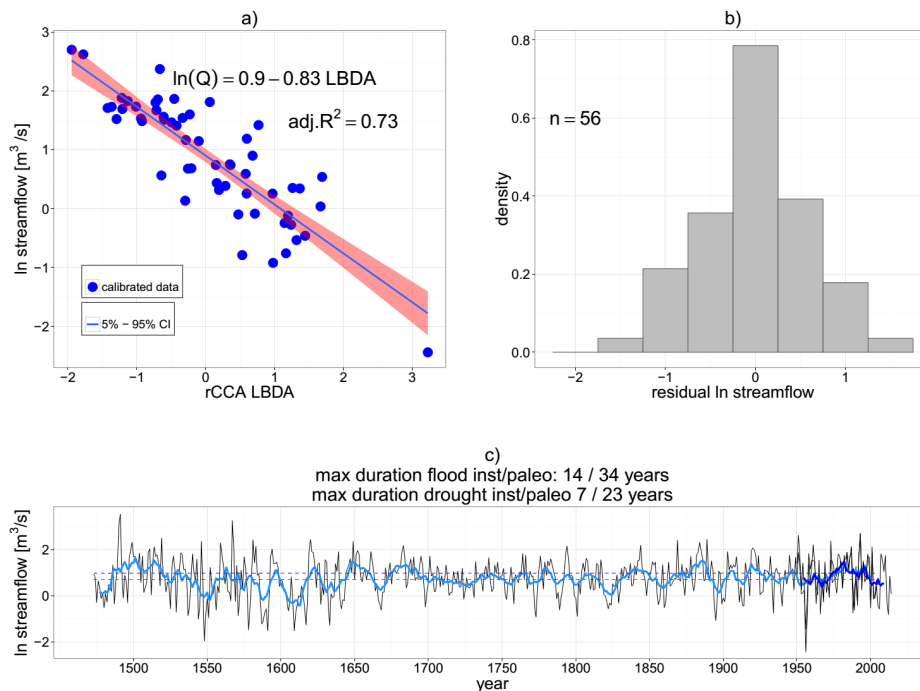
390 Where  $x_i$  and  $\hat{x}_i$  are the observed and modelled streamflows respectively, while  $\bar{x}_v$  and  $\bar{x}_c$  are the  
391 means of the observed streamflow in the validation and calibration periods respectively.

## 392 6 Results

393 A log-linear model was individually fitted to each streamflow gauge using logarithmically transformed  
394 streamflow and the first canonical variate of the LBDA inputs (schematic of inputs shown in Figure 3)  
395 using rCCA. The models were tested using a cross-validation procedure and calibration and verification  
396 metrics, as described in Section 5.3 were calculated. Finally, an overall assessment of the dominant  
397 modes of temporal variability in reconstructed streamflow variability in the Missouri River Basin was  
398 made using a frequency wavelet analysis on the leading PCs of reconstructed streamflow.

### 399 6.1 Model results

400 Streamflow at each of the 55 stations in the Missouri River Basin was constructed using a least squares  
401 regression model of natural-log streamflow and the first regularized canonical variate of the combined  
402 LBDA grids within a 450 km radius and first eight PCs of US-wide LBDA as the predictor variable. All  
403 results for modeled and reconstructed streamflow are shown in log space, while verification statistics  
404 are also calculated for logarithmic streamflow. An example of the input and modeled streamflow is  
405 shown for one streamflow station calibrated using data from Turkey Creek near Seneca (Figure 5). A  
406 summary of modeled streamflow results and summary statistics of the residuals for both calibrated and  
407 verified streamflow across all 100 calibration and verification sets are provided in the supporting  
408 information.



409

410 *Figure 5. Model results for Turkey Creek near Seneca (USGS gauge number 06814000) calibrated using all available data and the*  
 411 *first canonical variate of the LBDA inputs a) input (dots) and modelled (line) natural log streamflow showing 5<sup>th</sup> – 95<sup>th</sup> prediction*  
 412 *interval, b) histogram of model residuals, and c) flood and drought persistence gauged by a threshold of  $\text{mean} \pm 0.5SD$  of ten year*  
 413 *running average (blue lines – light blue line is the reconstructed ten year moving average).*

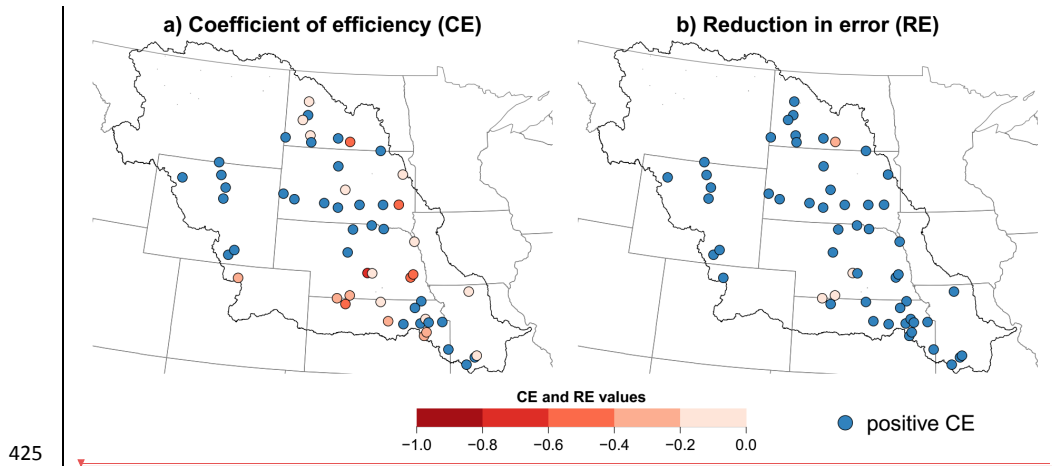
414 The rCCA linear model resulted in 55 models with adjusted  $R^2$  values ranging between 0.56 and 0.90.

415 The model residuals were near normal at almost all stations with the exception of some stations with  
 416 small catchment regions showing near-uniform residuals. In each streamflow model, the magnitude of  
 417 rCCA loadings for the eight PCs of CONUS LBDA were similar to those of the  $\sim 300$  local LBDA grids  
 418 indicating that information from the eight PCs of CONUS LBDA did not dominate the reconstruction.

## 419 6.2 Model validation

420 The predictive power of the model was assessed by calculating the coefficient of efficiency (CE) and  
 421 reduction of error (RE) for all cross validations. The median CE and RE values are shown in Figure 6 a and

422 b respectively, while distributions of the values are shown in box plots in the supporting information. CE  
423 and RE values range from  $-\infty$  to 1, with values over 0 indicating that the model predictions are more  
424 accurate than the respective climatology used for each statistic.

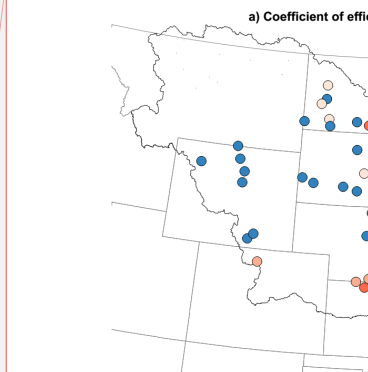


425  
426 *Figure 6. Median values of a) CE and b) RE values for cross validated models of the 55 Missouri River Basin streamflow stations.*

427 The CE values in Figure 6 indicate that many of the streamflow models are able to provide median  
428 streamflow predictions with greater accuracy than the climatology of the withheld data. The RE values  
429 likewise suggest that the models are able to predict the withheld values with greater accuracy than the  
430 climatology of the calibrated values at most stations. The combined CE and RE results suggest there is  
431 some skill in reconstructing paleoclimate streamflow using the LBDA.

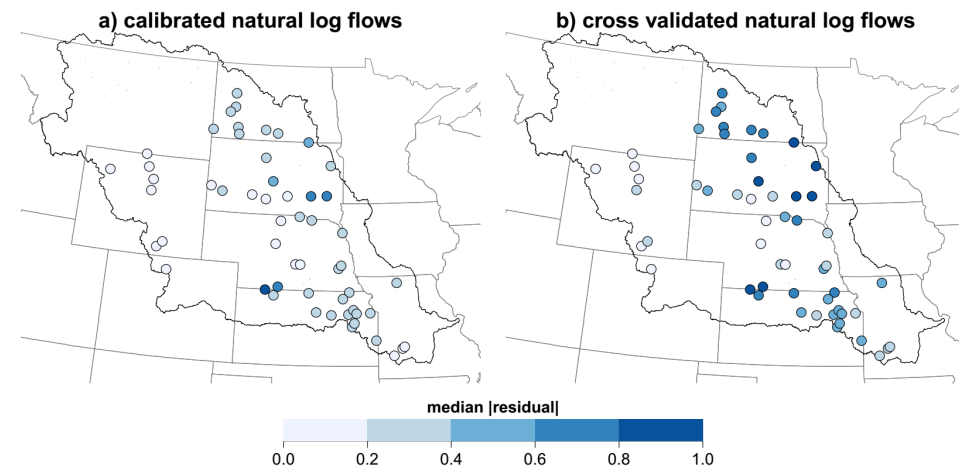
432 A comparison of the distribution of residuals resulting from calibrated values vs. validated values was  
433 also made using the 100 repetitions of the k-fold cross validation test. The median and the range of the  
434 cross-validated absolute residuals are modestly larger than that of the absolute residuals during the  
435 fitting process. The differences, accounting for the sample sizes, are not statistically significant in over  
436 90% of 100 cross validated results (using a two-sided t-test and a null hypothesis that the true difference  
437 between the means is zero and  $\alpha = 0.05$ ). The median values of the cross validated residuals are

Michelle Ho 3/6/2016 4:45 PM



Deleted:

439 shown in Figure 7, while box plots showing the distribution of the cross validated residuals are shown in  
440 the supporting information. A summary of the reconstruction is presented in Section 6.3 using PCs of the  
441 streamflow reconstructed at the 55 stations.

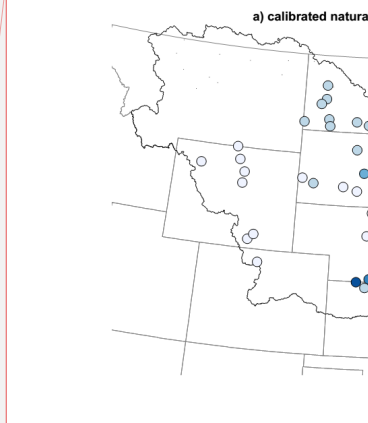


442  
443 *Figure 7. Median cross validated model residuals (absolute value) for a) calibrated and b) verified natural log streamflow inputs*  
444 *at the 55 Missouri River Basin streamflow stations. Cross validations were repeated 100 times.*

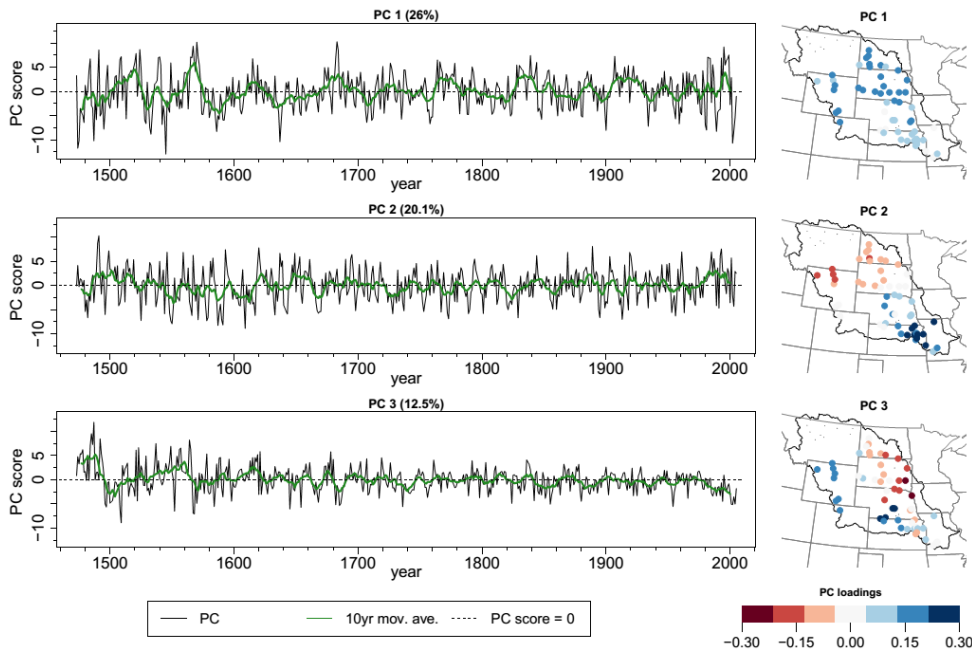
### 445 6.3 Reconstructed streamflow analysis

446 Mean annual streamflow was reconstructed for all 55 stations in the Missouri River Basin (data available  
447 in Ho et al., 2016). PCA was used to extract the leading modes of variability from the reconstructed  
448 Missouri River Basin natural log streamflow at all 55 stations. A combined assessment of the PCs using  
449 North's Rule of Thumb and a scree plot showed adjacent degenerate eigenvalues pairs for PCs of order 4  
450 and higher and discontinuities at PC4 respectively (shown in supporting information). Furthermore, the  
451 spatial pattern of PC4 was difficult to interpret and therefore only the first three PCs are shown in Figure  
452 8.

Michelle Ho 3/6/2016 4:46 PM



Deleted:



454

455 *Figure 8. Annual (black) and 10 year moving average (green line) of the first four PCs of reconstructed mean natural log*  
 456 *streamflow in the Missouri River Basin (left) with percentage variance explained shown in parenthesis and the corresponding*  
 457 *loading patterns (right).*

458 Negative streamflow anomalies coinciding with the 1930s dust bowl drought and 1950s drought are

459 evident in PC1, which has a positive loading pattern across all stations in the Missouri River Basin (Figure

460 8). The 1950s drought had more severe impacts in the southern half of the Missouri River Basin

461 (Piechota and Dracup, 1996; Cole et al., 2002; Andreadis et al., 2005; Cook et al., 2009) and this is

462 detected in PC2 that is comprised of positive loadings in the southern half of the Basin. The Civil War

463 drought, which largely impacted the Central Plains region from the mid-1850s to mid-1860s (Herweijer

464 et al., 2006), is also evident in PC2.

465 The severity of these droughts had wide ranging impacts on ecological states, agricultural produce and

466 social activities and are often used as benchmark droughts to which contemporary droughts are

467 compared (Breshears et al., 2005; Hornbeck, 2009). However, the streamflow reconstructions suggest

468 periods where streamflow deficits may have been more severe than any of these historical droughts. For  
469 example, the late 1540s, 1590s, and late 1750s all show negative streamflow anomalies in both PC1 and  
470 PC2 (26% and 20.1% variability explained respectively) that are of similar or greater magnitude than the  
471 streamflow deficits associated with the Civil War, Dust Bowl or 1950s droughts. Drought durations were  
472 also longer when assessed in the context of streamflow variability over the past 500 years. A threshold  
473 of above or below 0.5 standard deviation in decadal streamflow at each station was used as an  
474 approximation of flood and drought regimes respectively. The maximum duration of continuous periods  
475 of either flood or drought regimes were found to be longer in the majority of stations investigated as  
476 demonstrated in Figure 5 c that shows the reconstructed annual and decadal streamflow at Turkey  
477 Creek and the longest duration drought and flood regime for the station.

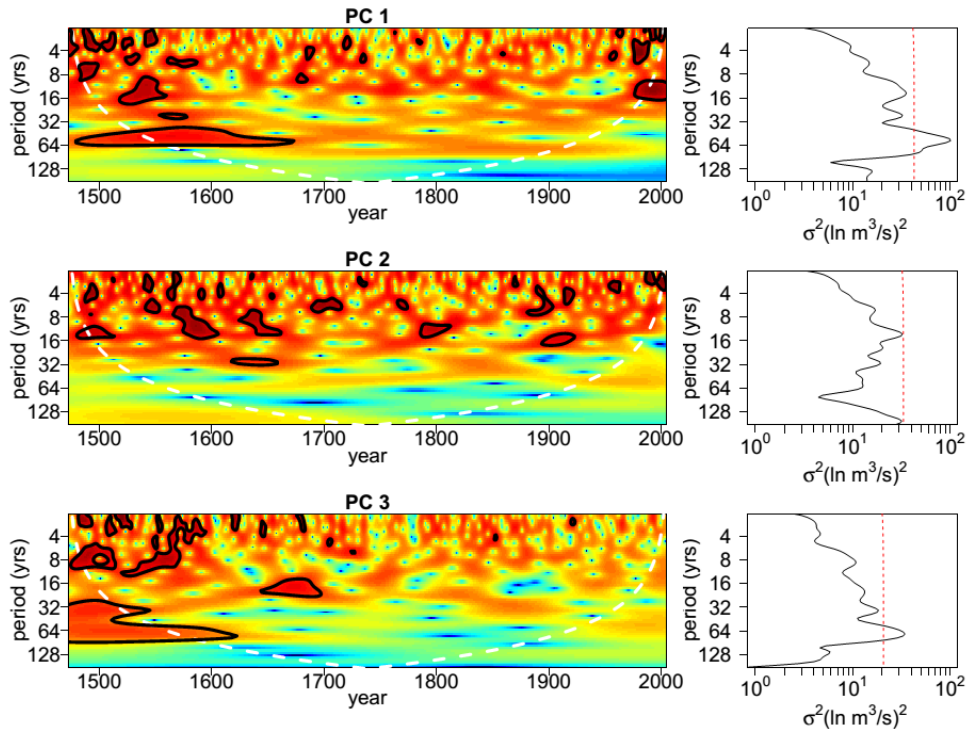
478 PC3 (12.5% variability explained) is a mode of variability largely influenced by differences in streamflow  
479 variability between the Rocky Mountains regions and the remainder of the Missouri River Basin. The  
480 apparent increase in variability in the 1500s and 1600s in PC3 suggests that differences in streamflow  
481 variability between the Rocky Mountains regions and the plains region may have been more  
482 pronounced during these two centuries.

483 PC2 of the reconstructed streamflow has positive loadings in the southern Missouri River Basin. This PC  
484 shows an extended period of below average flows, with the exception of a couple of years, in the first  
485 two decades of the 1600s and is evident in individual streamflow reconstructions around the border of  
486 Nebraska and South Dakota and in Kansas. A similar persistence of decadal-scale low flows are not  
487 featured in the instrumental record, suggesting that successive years of low streamflow have persisted  
488 for longer periods than what has been observed in historical records.

489 The PCs of the reconstructed streamflow also show several periods where annual flows are larger than  
490 the maximum instrumental streamflow. Periods of anomalously high streamflow include the 1560s for

491 PC 1 and the early 1490s for PC 2 and 3. These periods coincide with high positive reconstructed PDSI  
492 values particularly in the lower Missouri Basin region and along the middle and southern Rocky  
493 Mountains in the early 1490s and positive PDSI values across most of the interior plains. Anomalously  
494 large streamflow during the 1560s were reconstructed at stream gauges around the interior plains  
495 regions from tributaries that join the Missouri River in North and South Dakota (e.g. Cannonball, Grand  
496 and White Rivers). However, anomalously high streamflow around the 1560s was not reconstructed for  
497 either the Platte or Yellowstone rivers suggesting that the increase in streamflow may have resulted  
498 from weather systems delivering moisture in the upper Interior Plains (North and South Dakota), but not  
499 in the Rocky Mountains region.

500 In order to quantitatively assess the temporal modes of variability implied by the PCs, a continuous  
501 wavelet transform was applied to identify key periodicities in the main modes of reconstructed  
502 streamflow variability. Wavelet transforms for PCs 1-3 along with their global wavelet spectra are shown  
503 in Figure 9. The wavelet transforms were implemented using the Morlet wavelet function as the mother  
504 wavelet and padded with zeros to limit the edge effects of the wavelet analysis (Torrence and Compo,  
505 1998).



507

508 *Figure 9. Wavelet transforms of PC 1-3 of reconstructed log streamflow in the Missouri River Basin from 1473-2005 (left), black*  
 509 *lines show 95% significance level, dashed white line shows cone of influence under which edge effects may influence the results*  
 510 *and the time averaged global wavelet spectra (right) with red dashed line showing 95% confidence level for the wavelet*  
 511 *transform using a white noise background spectrum.*

512 All PCs show decadal-scale variability in the 1500s and 1600s, which is significant at the 95% level

513 against white noise for PC 1 and 3. Interestingly, none of the wavelet spectra show a consistent mode of

514 variability that coincide with an El Niño Southern Oscillation frequency (2-7 years Allan, 2000). In

515 addition, previously identified periods of decadal-scale variability in paleoclimate reconstructions of the

516 Pacific Decadal Oscillation (Gedalof and Smith, 2001; MacDonald and Case, 2005) are not reproduced in

517 the streamflow reconstruction. The multi-decadal-scale variability seen in PC1 and PC3 around the 16<sup>th</sup>

518 century suggest that the persistence of drought and pluvial events was longer around this period

519 compared with historic records and is reflected in previous reconstructions of persistent drought during

520 this period (Woodhouse and Overpeck, 1998) and in nearby regions (Meko et al., 1995). A similar  
521 pattern of multi-decadal-scale variability is also seen in the PCs of LBDA grids within the Missouri River  
522 Basin around the 16<sup>th</sup> century (results not shown here).

## 523 7 Discussion and Conclusions

524 Human activities and social processes are irrefutably intertwined with the hydrological cycle with each  
525 influencing and affecting the behavior of the other (Matalas et al., 1982). Consequently, hydrologic  
526 investigations need to occur at spatial scales that encompass the social and political regions influencing  
527 how water is stored, used, and transported. This requirement necessitates a broader spatial scale to be  
528 considered beyond the watershed scale typically addressed in hydrologic assessments (Vogel et al.,  
529 2015). Spatially and temporally broad measures of hydroclimatic variability such as the LBDA provide  
530 opportunities to conduct such assessments.

531 We used the LBDA to demonstrate a method of reconstructing mean annual streamflow in the Missouri  
532 River Basin from 1473-2005 in order to facilitate a future reconstruction of streamflow across the  
533 CONUS. The rCCA approach enabled a large degree of spatially coherent information in the LBDA to be  
534 condensed via linear transforms into a canonical variate that explained the majority of information in  
535 the target streamflow. The additional regularization step addressed the issue of an ill-posed problem  
536 where the number of variables far exceeded the number of observations. Regularization addressed  
537 potential issues with overfitting of the model and subsequent poor predictability of uncalibrated values.  
538 This site-by-site approach, opposed to a multi-site model, avoided data issues where streamflow records  
539 were not concurrently available at all sites. The method also tailored the LBDA information included in  
540 the model to each unique site and therefore has the potential to be applied to any individual streamflow  
541 site in the United States and, potentially, any streamflow site in North America where there are useful  
542 reconstructions in the LBDA.

543 The analysis of streamflow variability in the Missouri River Basin over the past 500 years was, however,  
544 limited by the absence of streamflow gauges in the headwaters of the Missouri River. Based on the  
545 selection criteria used here, there are no gauges located upstream of Fort Peck Lake and only one gauge  
546 in Montana was included in the analysis. The dearth of gauges in this region therefore omits potentially  
547 critical information when drawing conclusions regarding regional streamflow variability as presented in  
548 Section 6.3. Future analyses could seek to include additional streamflow gauges, potentially drawing on  
549 non-references gauges from the GAGESII database and other sources. While non-reference GAGESII  
550 gauges are located within disturbed catchments with potentially regulated streams, it may be plausible  
551 that these measurements would still preserve the signatures of interannual streamflow variability which  
552 we seek to recover. Encouragingly, the headwater regions of the Missouri River Basin contain a  
553 relatively dense network of tree-ring chronologies that would improve the fidelity of both PDSI and  
554 subsequent streamflow reconstructions. Furthermore, efforts towards expanding this tree-ring network  
555 are currently being made with a view to target and improve Missouri River Basin streamflow  
556 reconstructions (Pederson, 2013).

557 Dendrochronological studies for the purposes of climate reconstruction typically use a targeted  
558 sampling method whereby old living trees are selected for sampling to maximize the length of  
559 reconstruction and the location of sampling is chosen to ensure that tree growth is largely limited by the  
560 climatic factors of interest for reconstruction (Speer, 2010). Over 50 different tree species were used in  
561 the formulation of the LBDA, however, species was not a determinant for inclusion in the LBDA. Rather,  
562 the tree-ring chronologies specifically selected for use in reconstruction were those that were best  
563 correlated with available soil moisture as modeled a priori by PDSI and this was done using a method of  
564 point-by-point regression that ensures that the selected chronologies were likely to have a causal  
565 association with the gridded instrumental PDSI being reconstructed (Cook et al., 1999).

566 In developing our approach to reconstructing streamflow from the tree-ring based LBDA, we have  
567 capitalized on the fact that both variables are derivatives of a set of (unspecified) climate variables to  
568 which both streamflow and PDSI are sensitive (e.g. precipitation amount, temperature, and  
569 evaporation). Although variability is lost in each step of reconstructing PDSI from tree-rings and in  
570 reconstructing streamflow from the LBDA, it is possible that the reconstructed LBDA may remove noise  
571 irrelevant to hydrological variability rendering this paleoclimate reconstruction of PDSI particularly  
572 suitable for informing streamflow. Furthermore, streamflow and PDSI are inherently different  
573 hydrological variables as demonstrated in the temporal and spatial statistics shown in the supporting  
574 information. Streamflow is also a derivative of a suite of non-climatic variables such as catchment size,  
575 land-use, geology, topography, and groundwater interactions. These catchment-specific attributes are  
576 not explicitly considered in our study (with the exception of land use considered through the selection of  
577 streamflow gauges from relatively undisturbed catchments) and the inclusion of attributes could  
578 provide an improvement to the model (e.g. Thomas and Benson, 1970; Lima and Lall, 2010).

579 Assessments of hydrological variability on temporal scales broader than traditional hydrological  
580 practices are paramount to identifying long-term variations and quantifying the nature of persistent  
581 regimes, which are unachievable using comparatively short instrumental record. The understanding of  
582 paleoclimate variability facilitates the development of water and land use practices that are appropriate,  
583 sustainable, and resilient under previous patterns of hydroclimatic variability and complement efforts  
584 towards developing suitable adaption procedures for projected future climate scenarios. Whilst  
585 paleoclimate reconstructions of annual streamflow are useful for assessing long-term variability, annual  
586 information is often too coarse for water resource applications such as reservoir management and  
587 seasonal water allocations. In order to produce relevant data on a monthly or daily time step, a future  
588 undertaking could involve temporal disaggregation of the data. One method that could be utilized is a  
589 non-parametric k-nearest neighbors approach (Lall and Sharma, 1996; Rajagopalan and Lall, 1999;

590 Nowak et al., 2010) that would avoid some issues associated with multi-site stochastic generation  
591 schemes (Valencia and Schaake, 1973; Segond et al., 2006).The objective of producing paleoclimate  
592 reconstructions of streamflow that are relevant for water resource management also suggests that the  
593 reconstruction of streamflow that is currently regulated would be valuable. Paleoclimate  
594 reconstructions of regulated streamflow could be achieved at an annual scale provided that regulations  
595 primarily impact the timing of sub-annual flows or if sufficient data prior to regularization exists. The  
596 resulting annual streamflow reconstruction could be disaggregated to a shorter timescale if pre-  
597 regulated flow data is available.

598 This study has capitalized on the availability of the spatially and temporally complete LBDA dataset to  
599 attempt a reconstruction of streamflow on a broad spatial scale across the Missouri River Basin. The  
600 analysis employed here provides a feasible foundation from which streamflow reconstructions across  
601 the CONUS and North America could be implemented using the LBDA. In addition, the methodology  
602 presented here could be applied to reconstructing streamflow using other reconstructions of drought  
603 including the Monsoon Asia Drought Atlas (Cook et al., 2010b), the Old World Drought Atlas covering  
604 Europe, North Africa, and the Middle East (Cook et al., 2015b) and the Australia and New Zealand  
605 summer drought atlas (Palmer et al., 2015).

606

Michelle Ho 3/6/2016 6:27 PM

**Deleted:** similar

Michelle Ho 3/6/2016 6:28 PM

**Deleted:** reconstructions could be attempted

Michelle Ho 3/6/2016 6:28 PM

**Deleted:** continental-scale

Michelle Ho 3/6/2016 6:28 PM

**Deleted:** s

Michelle Ho 3/6/2016 6:29 PM

**Deleted:** and

Michelle Ho 3/6/2016 6:29 PM

**Deleted:** recently completed

## 614 Acknowledgments

615 We would like to sincerely thank Ben Cook (bc9z@ldeo.columbia.edu) for providing the LBDA data, the  
616 US Geological Survey (USGS) for making their monthly streamflow data available  
617 ([http://waterdata.usgs.gov/nwis/monthly?referred\\_module=sw](http://waterdata.usgs.gov/nwis/monthly?referred_module=sw)), staff at the USGS (Christopher Ryan  
618 cmryan@usgs.gov and Thomas Weaver tlweaver@usgs.gov) for clarifying streamflow data availability;  
619 Siyan Wang for initial data collection; ~~Scott Steinschneider, Xun Sun and Pierre Gentine at the Columbia~~  
620 ~~Water Center for their invaluable discussions on this work; and Juan A. Ballesteros and one other~~  
621 ~~anonymous reviewer for their constructive comments.~~ Reconstruction results are available on the NOAA  
622 paleoclimate data base (<https://www.ncdc.noaa.gov/paleo/study/19520>). This work is funded by an NSF  
623 award 1360446 and NSF award 1401698. Lamont-Doherty contribution number XXXX.

## 624 References

- 625 Allan, R. J., 2000: ENSO and Climatic Variability in the Past 150 Years. *El Niño and the Southern*  
626 *Oscillation: Multiscale Variability and Global and Regional Impacts*, Henry F. Diaz and Vera  
627 Markgraf, Eds., Cambridge University Press, pp.  
628 Allen, E. B., Rittenour, T. M., DeRose, R. J., Bekker, M. F., Kjelgren, R. and Buckley, B. M., 2013: A tree-  
629 ring based reconstruction of Logan River streamflow, northern Utah, *Water Resources Research*,  
630 **49** (12), 8579-8588, doi: 10.1002/2013wr014273.  
631 Andreadis, K. M., Clark, E. A., Wood, A. W., Hamlet, A. F. and Lettenmaier, D. P., 2005: Twentieth-  
632 Century Drought in the Conterminous United States, *Journal of Hydrometeorology*, **6** (6), 985-  
633 1001, doi: 10.1175/JHM450.1.  
634 Ault, T. R., Cole, J. E., Overpeck, J. T., Pederson, G. T. and Meko, D. M., 2014: Assessing the risk of  
635 persistent drought using climate model simulations and paleoclimate data, *Journal of Climate*,  
636 **27** (20), 7529-7549, doi: 10.1175/jcli-d-12-00282.1.  
637 Breshears, D. D., Cobb, N. S., Rich, P. M., Price, K. P., Allen, C. D., Balice, R. G., Romme, W. H., Kastens, J.  
638 H., Floyd, M. L., Belnap, J., Anderson, J. J., Myers, O. B. and Meyer, C. W., 2005: Regional  
639 vegetation die-off in response to global-change-type drought, *Proceedings of the National*  
640 *Academy of Sciences of the United States of America*, **102** (42), 15144-15148, doi:  
641 10.1073/pnas.0505734102.  
642 Brown, D. P. and Comrie, A. C., 2004: A winter precipitation 'dipole' in the western United States  
643 associated with multidecadal ENSO variability, *Geophysical Research Letters*, **31** (9), doi:  
644 10.1029/2003GL018726.  
645 Buuren, S. v. and Groothuis-Oudshoorn, K., 2011: mice: Multivariate Imputation by Chained Equations in  
646 R, *Journal of Statistical Software*, **45** (3), 67.

Michelle Ho 4/6/2016 7:48 AM

Deleted: and

648 Cole, J. E., Overpeck, J. T. and Cook, E. R., 2002: Multiyear La Niña events and persistent drought in the  
649 contiguous United States, *Geophysical Research Letters*, **29** (13), 25-1-25-4, doi:  
650 10.1029/2001GL013561.

651 Cook, B. I., Miller, R. L. and Seager, R., 2009: Amplification of the North American “Dust Bowl” drought  
652 through human-induced land degradation, *Proceedings of the National Academy of Sciences*,  
653 **106** (13), 4997-5001, doi: 10.1073/pnas.0810200106.

654 Cook, B. I., Ault, T. R. and Smerdon, J. E., 2015a: Unprecedented 21st century drought risk in the  
655 American Southwest and Central Plains, *Science Advances*, **1** (1), doi: 10.1126/sciadv.1400082.

656 Cook, B. I., Smerdon, J. E., Seager, R. and Cook, E. R., 2013a: Pan-Continental Droughts in North America  
657 over the Last Millennium, *Journal of Climate*, **27** (1), 383-397, doi: 10.1175/jcli-d-13-00100.1.

658 Cook, E. R. and Krusic, P. J., 2004: The North American Drought Atlas,  
659 Cook, E. R., Briffa, K. R. and Jones, P. D., 1994: Spatial regression methods in dendroclimatology: A  
660 review and comparison of two techniques, *International Journal of Climatology*, **14** (4), 379-402,  
661 doi: 10.1002/joc.3370140404.

662 Cook, E. R., Meko, D. M., Stahle, D. W. and Cleaveland, M. K., 1999: Drought Reconstructions for the  
663 Continental United States, *Journal of Climate*, **12** (4), 1145-1162, doi: 10.1175/1520-  
664 0442(1999)012<1145:drftcu>2.0.co;2.

665 Cook, E. R., Woodhouse, C. A., Eakin, C. M., Meko, D. M. and Stahle, D. W., 2004: Long-Term Aridity  
666 Changes in the Western United States, *Science*, **306** (5698), 1015-1018, doi:  
667 10.1126/science.1102586.

668 Cook, E. R., Seager, R., Heim, R. R., Vose, R. S., Herweijer, C. and Woodhouse, C., 2010a: Megadroughts  
669 in North America: placing IPCC projections of hydroclimatic change in a long-term palaeoclimate  
670 context, *Journal of Quaternary Science*, **25** (1), 48-61, doi: 10.1002/jqs.1303.

671 Cook, E. R., Anchukaitis, K. J., Buckley, B. M., D’Arrigo, R. D., Jacoby, G. C. and Wright, W. E., 2010b:  
672 Asian Monsoon Failure and Megadrought During the Last Millennium, *Science*, **328** (5977), 486-  
673 489, doi: 10.1126/science.1185188.

674 Cook, E. R., Palmer, J. G., Ahmed, M., Woodhouse, C. A., Fenwick, P., Zafar, M. U., Wahab, M. and Khan,  
675 N., 2013b: Five centuries of Upper Indus River flow from tree rings, *Journal of Hydrology*, **486**  
676 (0), 365-375, doi: 10.1016/j.jhydrol.2013.02.004.

677 Cook, E. R., Seager, R., Kushnir, Y., Briffa, K. R., Büntgen, U., Frank, D., Krusic, P. J., Tegel, W., van der  
678 Schrier, G., Andreu-Hayles, L., Baillie, M., Baittinger, C., Bleicher, N., Bonde, N., Brown, D.,  
679 Carrer, M., Cooper, R., Čufar, K., Dittmar, C., Esper, J., Griggs, C., Gunnarson, B., Günther, B.,  
680 Gutierrez, E., Haneca, K., Helama, S., Herzig, F., Heussner, K.-U., Hofmann, J., Janda, P., Kontic,  
681 R., Köse, N., Kyncl, T., Levanič, T., Linderholm, H., Manning, S., Melvin, T. M., Miles, D., Neuwirth,  
682 B., Nicolussi, K., Nola, P., Panayotov, M., Popa, I., Rothe, A., Seftigen, K., Seim, A., Svarva, H.,  
683 Svoboda, M., Thun, T., Timonen, M., Touchan, R., Trotsiuk, V., Trouet, V., Walder, F., Ważny, T.,  
684 Wilson, R. and Zang, C., 2015b: Old World megadroughts and pluvials during the Common Era,  
685 *Science Advances*, **1** (10), doi: 10.1126/sciadv.1500561.

686 Cutler, A. and Breiman, L., 1994: Archetypal Analysis, *Technometrics*, **36** (4), 338-347, doi:  
687 10.1080/00401706.1994.10485840.

688 Dai, A., 2013: Increasing drought under global warming in observations and models, *Nature Clim.*  
689 *Change*, **3** (1), 52-58, doi: 10.1038/nclimate1633.

690 De Bie, T. and De Moor, B., 2003: On the regularization of canonical correlation analysis, *Int. Sympos. ICA*  
691 *and BSS*, 785-790.

692 Dettinger, M. D., Cayan, D. R., Diaz, H. F. and Meko, D. M., 1998: North–South Precipitation Patterns in  
693 Western North America on Interannual-to-Decadal Timescales, *Journal of Climate*, **11** (12), 3095-  
694 3111, doi: 10.1175/1520-0442(1998)011<3095:NSPPIW>2.0.CO;2.

695 Devineni, N., Lall, U., Pederson, N. and Cook, E., 2013: A Tree-Ring-Based Reconstruction of Delaware  
696 River Basin Streamflow Using Hierarchical Bayesian Regression, *Journal of Climate*, **26** (12),  
697 4357-4374, doi: 10.1175/jcli-d-11-00675.1.

698 Dijkstra, T., 2014: Ridge regression and its degrees of freedom, *Quality & Quantity*, **48** (6), 3185-3193,  
699 doi: 10.1007/s11135-013-9949-7.

700 Driscoll, D., 2013: *Hydroclimatic information and risk assessment methods for managing hydrologic*  
701 *extremes in the Missouri River Basin - A white paper for future studies. Hydroclimatic extremes in*  
702 *the Missouri River Basin*, USGS.

703 Fritts, H. C., 1976: *Tree Rings and Climate*. Academic Press Inc., 567 pp.

704 Galat, D. L., Berry Jr, C. R., Peters, E. J. and White, R. G., 2005: Chapter 10 - Missouri River Basin. *Rivers of*  
705 *North America*, Arthur C. Benke and Colbert E. Cushing, Eds., Academic Press, pp. 426-480.

706 Gallant, A. J. E. and Gergis, J., 2011: An experimental streamflow reconstruction for the River Murray,  
707 Australia, 1783-1988, *Water Resources Research*, **47** (12), W00G04, doi:  
708 10.1029/2010wr009832.

709 Gedalof, Z. e. and Smith, D. J., 2001: Interdecadal climate variability and regime-scale shifts in Pacific  
710 North America, *Geophys. Res. Lett.*, **28** (8), 1515-1518, doi: 10.1029/2000gl011779.

711 Gershunov, A. and Barnett, T. P., 1998: Interdecadal Modulation of ENSO Teleconnections, *Bulletin of*  
712 *the American Meteorological Society*, **79** (12), 2715-2725, doi: 10.1175/1520-  
713 0477(1998)079<2715:IMOET>2.0.CO;2.

714 González, I., Déjean, S., Martin, P. G. P. and Baccini, A., 2008: CCA: An R Package to Extend Canonical  
715 Correlation Analysis, *2008*, **23** (12), 14, doi: 10.18637/jss.v023.i12.

716 Herweijer, C., Seager, R. and Cook, E. R., 2006: North American droughts of the mid to late nineteenth  
717 century: a history, simulation and implication for Mediaeval drought, *The Holocene*, **16** (2), 159-  
718 171, doi: 10.1191/0959683606hl917rp.

719 Herweijer, C., Seager, R., Cook, E. R. and Emile-Geay, J., 2007: North American Droughts of the Last  
720 Millennium from a Gridded Network of Tree-Ring Data, *Journal of Climate*, **20** (7), 1353-1376,  
721 doi: 10.1175/JCLI4042.1.

722 Higgins, R. W., Yao, Y., Yarosh, E. S., Janowiak, J. E. and Mo, K. C., 1997: Influence of the Great Plains  
723 Low-Level Jet on Summertime Precipitation and Moisture Transport over the Central United  
724 States, *Journal of Climate*, **10** (3), 481-507, doi: 10.1175/1520-  
725 0442(1997)010<0481:IOTGPL>2.0.CO;2.

726 Ho, M., Kiem, A. S. and Verdon, D. C., 2015: A paleoclimate rainfall reconstruction in the Murray-Darling  
727 Basin (MDB), Australia: 2. Assessing hydroclimatic risk using paleoclimate records of wet and dry  
728 epochs, *Water Resour. Res.*, **51** (10), doi: 10.1002/2015WR017059.

729 Ho, M., Lall, U. and Cook, E. R., 2016: Missouri River Basin 533 Year Annual Streamflow  
730 Reconstructions, <<https://www.ncdc.noaa.gov/paleo/study/19520>>.

731 Hornbeck, R., 2009: The Enduring Impact of the American Dust Bowl: Short and Long-run Adjustments to  
732 Environmental Catastrophe, *National Bureau of Economic Research Working Paper Series*, **No.**  
733 **15605**, doi: 10.3386/w15605.

734 Hotelling, H., 1936: Relations between two sets of variates, *Biometrika*, **28** (3-4), 321-377, doi:  
735 10.1093/biomet/28.3-4.321.

736 Jolliffe, I. T., 2002: *Principal Component Analysis*. Second edition ed. Springer-Verlag, 487 pp.

737 Jones, P. D. and Mann, M. E., 2004: Climate over past millennia, *Rev. Geophys.*, **42** (2), RG2002, doi:  
738 10.1029/2003rg000143.

739 Kunkel, K. E., Stevens, L. E., Stevens, S. E., Sun, L., Janssen, E., Wuebbles, D., Kruk, M. C., Thomas, D. P.,  
740 Shulski, M., Umphlett, N., Hubbard, K., Robbins, K., Romolo, L., Akyuz, A., Pathak, T., Bergantino,  
741 T. and Dobson, J. G., 2013: Regional Climate Trends and Scenarios for the U.S. National Climate  
742 Assessment. Part 4. Climate of the U.S. Great Plains, NOAA Technical Report, NOAA, pp. 82.

743 Lall, U. and Sharma, A., 1996: A Nearest Neighbor Bootstrap For Resampling Hydrologic Time Series,  
744 *Water Resources Research*, **32** (3), 679-693, doi: 10.1029/95WR02966.

745 Leurgans, S. E., Moyeed, R. A. and Silverman, B. W., 1993: Canonical Correlation Analysis when the Data  
746 are Curves, *Journal of the Royal Statistical Society. Series B (Methodological)*, **55** (3), 725-740.

747 Lima, C. H. R. and Lall, U., 2010: Spatial scaling in a changing climate: A hierarchical bayesian model for  
748 non-stationary multi-site annual maximum and monthly streamflow, *Journal of Hydrology*, **383**  
749 (3-4), 307-318, doi: 10.1016/j.jhydrol.2009.12.045.

750 Lockwood, J. G., 1999: Is Potential Evapotranspiration and Its Relationship with Actual  
751 Evapotranspiration Sensitive to Elevated Atmospheric CO<sub>2</sub> Levels?, *Climatic Change*, **41** (2), 193-  
752 212, doi: 10.1023/A:1005469416067.

753 MacDonald, G. M. and Case, R. A., 2005: Variations in the Pacific Decadal Oscillation over the past  
754 millennium, *Geophys. Res. Lett.*, **32** (8), L08703, doi: 10.1029/2005gl022478.

755 Matalas, N. C., Landwehr, J. M. and Wolman, M. G., 1982: Prediction in Water Management. *Scientific*  
756 *basis of water management*, National Academy, pp. 118-122.

757 McGowan, H. A., Marx, S. K., Denholm, J., Soderholm, J. and Kamber, B. S., 2009: Reconstructing annual  
758 inflows to the headwater catchments of the Murray River, Australia, using the Pacific Decadal  
759 Oscillation, *Geophys. Res. Lett.*, **36** (6), L06707, doi: 10.1029/2008gl037049.

760 Meko, D., Stockton, C. W. and Boggess, W. R., 1995: The tree-ring record of sever sustained drought,  
761 *JAWRA Journal of the American Water Resources Association*, **31** (5), 789-801, doi:  
762 10.1111/j.1752-1688.1995.tb03401.x.

763 Monteith, J., 1965: Evaporation and environment, *Symp. Soc. Exp. Biol.*, **4**.

764 Nohara, D., Kitoh, A., Hosaka, M. and Oki, T., 2006: Impact of Climate Change on River Discharge  
765 Projected by Multimodel Ensemble, *Journal of Hydrometeorology*, **7** (5), 1076-1089, doi:  
766 10.1175/JHM531.1.

767 Nowak, K., Prairie, J., Rajagopalan, B. and Lall, U., 2010: A nonparametric stochastic approach for  
768 multisite disaggregation of annual to daily streamflow, *Water Resources Research*, **46** (8),  
769 W08529, doi: 10.1029/2009WR008530.

770 Palmer, J. G., Cook, E. R., Turney, C. S. M., Allen, K., Fenwick, P., Cook, B. I., O'Donnell, A., Lough, J.,  
771 Grierson, P. and Baker, P., 2015: Drought variability in the eastern Australia and New Zealand  
772 summer drought atlas (ANZDA, CE 1500–2012) modulated by the Interdecadal Pacific  
773 Oscillation, *Environmental Research Letters*, **10** (12), 124002, doi: 10.1088/1748-  
774 9326/10/12/124002.

775 Pederson, G. T., 2013, Multi-century perspectives on current and future streamflow in the Missouri River  
776 Basin, viewed 2015/09/08, <

777 Pederson, N., Bell, A. R., Cook, E. R., Lall, U., Devineni, N., Seager, R., Eggleston, K. and Vranes, K. P.,  
778 2012: Is an Epic Pluvial Masking the Water Insecurity of the Greater New York City Region?,  
779 *Journal of Climate*, **26** (4), 1339-1354, doi: 10.1175/jcli-d-11-00723.1.

780 Piechota, T. C. and Dracup, J. A., 1996: Drought and Regional Hydrologic Variation in the United States:  
781 Associations with the El Niño-Southern Oscillation, *Water Resources Research*, **32** (5), 1359-  
782 1373, doi: 10.1029/96WR00353.

783 Pizarro, G. and Lall, U., 2002: El Niño-induced Flooding in the U.S. West: What Can We Expect, *Eos Trans.*  
784 *AGU*, **82** (32), 349-352, doi: 10.1029/2002EO000255.

785 Prairie, J., Nowak, K., Rajagopalan, B., Lall, U. and Fulp, T., 2008: A stochastic nonparametric approach  
786 for streamflow generation combining observational and paleoreconstructed data, *Water*  
787 *Resources Research*, **44** (6), W06423, doi: 10.1029/2007wr006684.

788 Quinn, W. H., 1992: A study of the Southern Oscillation-related climatic activity for A.D 622-1990  
789 incorporating Nile River flood data. *El Niño: Historical and Paleoclimatic Aspects of the Southern*  
790 *Ocean*, Henry F. Diaz and Vera Markgraf, Eds., Cambridge University Press, pp. 119-149.

791 Rajagopalan, B. and Lall, U., 1999: A k-nearest-neighbor simulator for daily precipitation and other  
792 weather variables, *Water Resources Research*, **35** (10), 3089-3101, doi:  
793 10.1029/1999WR900028.

794 Redmond, K. T. and Koch, R. W., 1991: Surface Climate and Streamflow Variability in the Western United  
795 States and Their Relationship to Large-Scale Circulation Indices, *Water Resources Research*, **27**  
796 (9), 2381-2399, doi: 10.1029/91WR00690.

797 Routson, C. C., Woodhouse, C. A. and Overpeck, J. T., 2011: Second century megadrought in the Rio  
798 Grande headwaters, Colorado: How unusual was medieval drought?, *Geophys. Res. Lett.*, **38**  
799 (22), L22703, doi: 10.1029/2011gl050015.

800 Seager, R., Ting, M., Held, I., Kushnir, Y., Lu, J., Vecchi, G., Huang, H.-P., Harnik, N., Leetmaa, A., Lau, N.-  
801 C., Li, C., Velez, J. and Naik, N., 2007: Model Projections of an Imminent Transition to a More  
802 Arid Climate in Southwestern North America, *Science*, **316** (5828), 1181-1184, doi:  
803 10.1126/science.1139601.

804 Segond, M. L., Onof, C. and Wheeler, H. S., 2006: Spatial-temporal disaggregation of daily rainfall from a  
805 generalized linear model, *Journal of Hydrology*, **331** (3-4), 674-689, doi:  
806 10.1016/j.jhydrol.2006.06.019.

807 Sheffield, J., Wood, E. F. and Roderick, M. L., 2012: Little change in global drought over the past 60  
808 years, *Nature*, **491** (7424), 435-438, doi: 10.1038/nature11575.

809 Smerdon, J. E., Cook, B. I., Cook, E. R. and Seager, R., 2015: Bridging Past and Future Climate across  
810 Paleoclimatic Reconstructions, Observations, and Models: A Hydroclimate Case Study, *Journal of*  
811 *Climate*, **28** (8), 3212-3231, doi: 10.1175/JCLI-D-14-00417.1.

812 Smith, S. R., Green, P. M., Leonardi, A. P. and O'Brien, J. J., 1998: Role of Multiple-Level Tropospheric  
813 Circulations in Forcing ENSO Winter Precipitation Anomalies, *Monthly Weather Review*, **126**  
814 (12), 3102-3116, doi: 10.1175/1520-0493(1998)126<3102:ROMLTC>2.0.CO;2.

815 Speer, J. H., 2010: *Fundamentals of tree-ring research*. University of Arizona Press, 333 pp.

816 St. George, S., 2010: Tree Rings as Paleoflood and Paleostage Indicators. *Tree Rings and Natural*  
817 *Hazards: A State-of-Art*, Markus Stoffel, Michelle Bollschweiler, R. David Butler, and H. Brian  
818 Luckman, Eds., Springer Netherlands, pp. 233-239.

819 Steinschneider, S. and Lall, U., 2015: Daily Precipitation and Tropical Moisture Exports across the Eastern  
820 United States: An Application of Archetypal Analysis to Identify Spatiotemporal Structure,  
821 *Journal of Climate*, **28** (21), 8585-8602, doi: 10.1175/JCLI-D-15-0340.1.

822 Stone, E. and Cutler, A., 1996: Introduction to archetypal analysis of spatio-temporal dynamics, *Physica*  
823 *D: Nonlinear Phenomena*, **96** (1-4), 110-131, doi: 10.1016/0167-2789(96)00016-4.

824 Thomas, D. M. and Benson, M. A., 1970: Generalization of streamflow characteristics from drainage-  
825 basin characteristics, Water Supply Paper, 1975, pp.

826 Thornthwaite, C. W., 1948: An Approach toward a Rational Classification of Climate, *Geographical*  
827 *Review*, **38** (1), 55-94, doi: 10.2307/210739.

828 Tierney, J. E., Oppo, D. W., Rosenthal, Y., Russell, J. M. and Linsley, B. K., 2010: Coordinated hydrological  
829 regimes in the Indo - Pacific region during the past two millennia, *Paleoceanography*, **25**,  
830 PA1102, doi: 10.1029/2009PA001871.

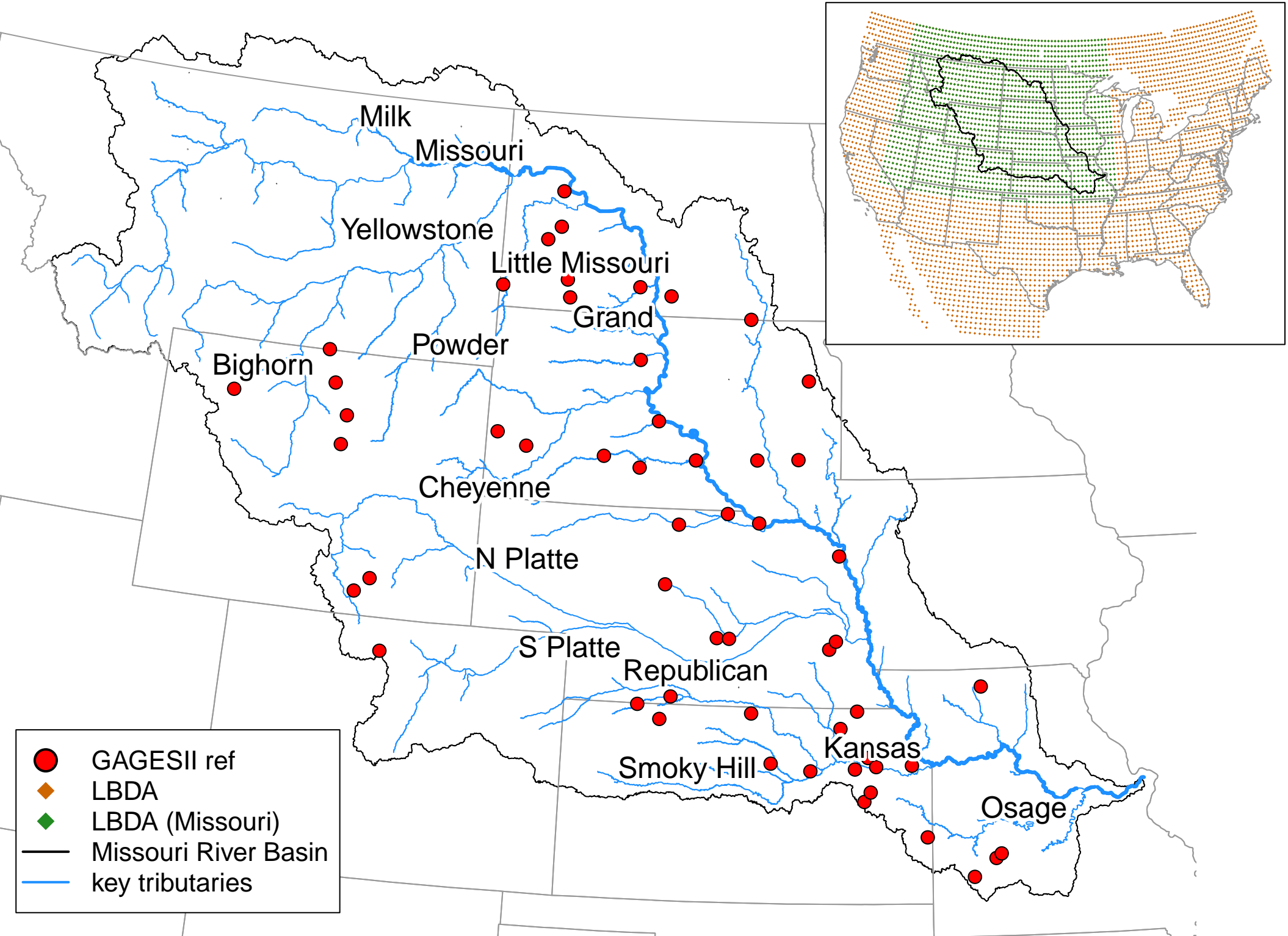
831 Tootle, G. A. and Piechota, T. C., 2006: Relationships between Pacific and Atlantic ocean sea surface  
832 temperatures and U.S. streamflow variability, *Water Resources Research*, **42** (7), n/a-n/a, doi:  
833 10.1029/2005WR004184.

834 Torrence, C. and Compo, G. P., 1998: A Practical Guide to Wavelet Analysis, *Bulletin of the American*  
835 *Meteorological Society*, **79** (1), 61-78, doi: 10.1175/1520-  
836 0477(1998)079<0061:APGTWA>2.0.CO;2.

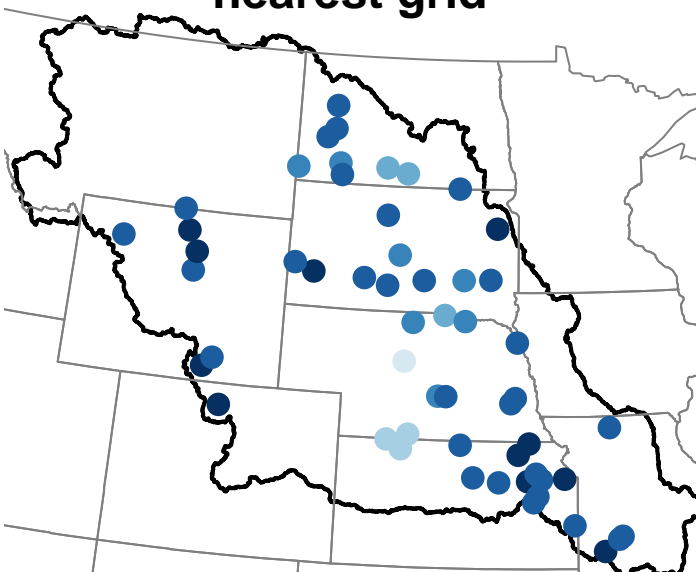
837 Trenberth, K. E., Branstator, G. W. and Arkin, P. A., 1988: Origins of the 1988 North American Drought,  
838 *Science*, **242** (4886), 1640-1645, doi: 10.1126/science.242.4886.1640.

839 U.S. Geological Survey, 2011, GAGES-II: Geospatial Attributes of Gages for Evaluating Streamflow,  
840 viewed 06/19/2015,  
841 <[http://water.usgs.gov/GIS/metadata/usgswrd/XML/gagesII\\_Sept2011.xml](http://water.usgs.gov/GIS/metadata/usgswrd/XML/gagesII_Sept2011.xml)>.  
842 US Army Corps of Engineers, 2006: Missouri River Mainstem Reservoir System - Master Water Control  
843 Manual - Missouri River Basin, pp. 431.  
844 Valencia, D. and Schaake, J. C., 1973: Disaggregation processes in stochastic hydrology, *Water Resour.*  
845 *Res*, **9** (3), 580-585, doi: 10.1029/WR009i003p00580.  
846 Vance, T. R., van Ommen, T. D., Curran, M. A. J., Plummer, C. T. and Moy, A. D., 2012: A Millennial Proxy  
847 Record of ENSO and Eastern Australian Rainfall from the Law Dome Ice Core, East Antarctica,  
848 *Journal of Climate*, **26** (3), 710-725, doi: 10.1175/jcli-d-12-00003.1.  
849 Vinod, H. D., 1976: Canonical ridge and econometrics of joint production, *Journal of Econometrics*, **4** (2),  
850 147-166, doi: 10.1016/0304-4076(76)90010-5.  
851 Vogel, R. M., Lall, U., Cai, X., Rajagopalan, B., Weiskel, P., Hooper, R. P. and Matalas, N. C., 2015:  
852 Hydrology: The interdisciplinary science of water, *Water Resources Research*, **51** (6), 4409-4430,  
853 doi: 10.1002/2015WR017049.  
854 White, I. R., Royston, P. and Wood, A. M., 2011: Multiple imputation using chained equations: Issues and  
855 guidance for practice, *Statistics in Medicine*, **30** (4), 377-399, doi: 10.1002/sim.4067.  
856 Wilson, R., Cook, E., D'Arrigo, R., Riedwyl, N., Evans, M. N., Tudhope, A. and Allan, R., 2010:  
857 Reconstructing ENSO: the influence of method, proxy data, climate forcing and teleconnections,  
858 *Journal of Quaternary Science*, **25** (1), 62-78, doi: 10.1002/jqs.1297.  
859 Wise, E. K., 2010: Spatiotemporal variability of the precipitation dipole transition zone in the western  
860 United States, *Geophysical Research Letters*, **37** (7), doi: 10.1029/2009GL042193.  
861 Woodhouse, C. A. and Meko, D., 1997: Number of Winter Precipitation Days Reconstructed from  
862 Southwestern Tree Rings, *Journal of Climate*, **10** (10), 2663-2669, doi: 10.1175/1520-  
863 0442(1997)010<2663:NOWPDR>2.0.CO;2.  
864 Woodhouse, C. A. and Overpeck, J. T., 1998: 2000 years of drought variability in the central United  
865 States, *Bulletin of the American Meteorological Society*, **79** (12), 2693, doi.  
866 Woodhouse, C. A., Gray, S. T. and Meko, D. M., 2006: Updated streamflow reconstructions for the Upper  
867 Colorado River Basin, *Water Resources Research*, **42** (5), W05415, doi: 10.1029/2005wr004455.  
868 Woodhouse, C. A., Kunkel, K. E., Easterling, D. R. and Cook, E. R., 2005: The twentieth-century pluvial in  
869 the western United States, *Geophysical Research Letters*, **32** (7), doi: 10.1029/2005GL022413.  
870 Woodhouse, C. A., Meko, D. M., MacDonald, G. M., Stahle, D. W. and Cook, E. R., 2010: A 1,200-year  
871 perspective of 21st century drought in southwestern North America, *Proceedings of the National*  
872 *Academy of Sciences*, **107** (50), 21283-21288, doi: 10.1073/pnas.0911197107.  
873 Woodhouse, C. A., Lukas, J., Brice, B., Hirschboeck, K., Hartmann, H., Kostuk, M., Lay, E., Martinex, D.,  
874 McMahan, B. and Shah, A., 2002, Tree rings and streamflow, viewed 2015/06/01,  
875 <<http://treeflow.info/content/tree-rings-and-streamflow>>.

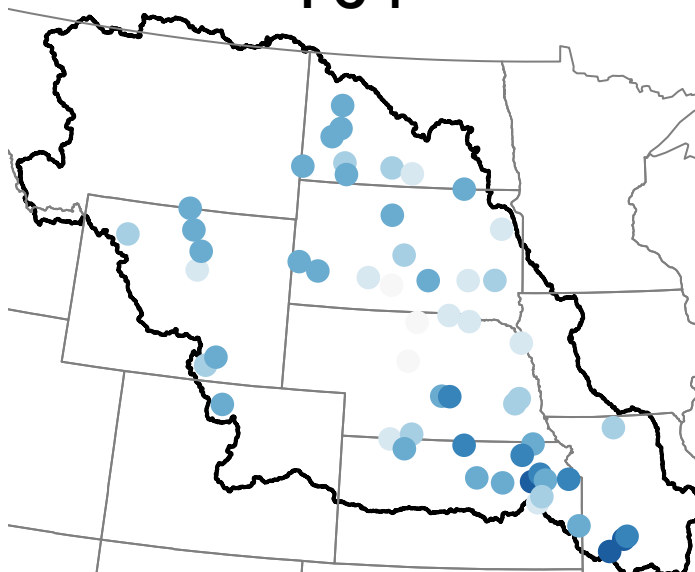
876



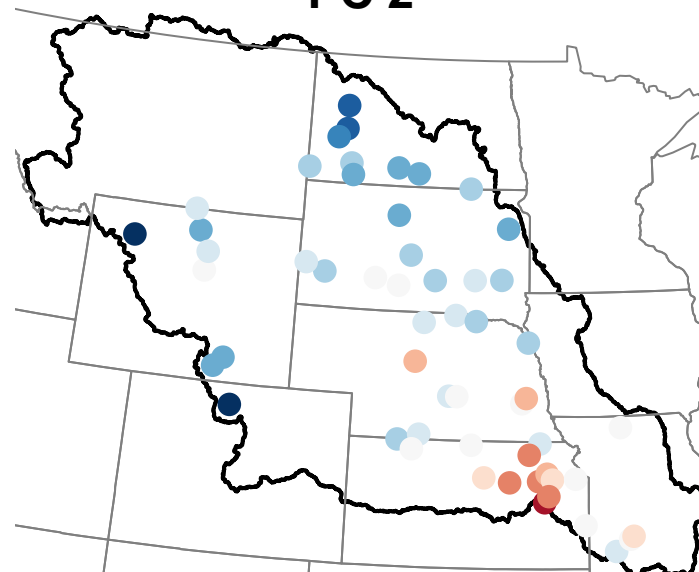
nearest grid



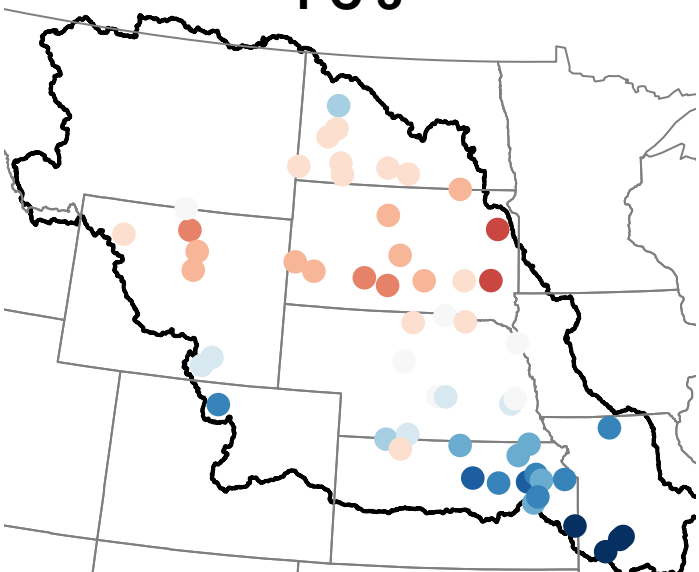
PC 1



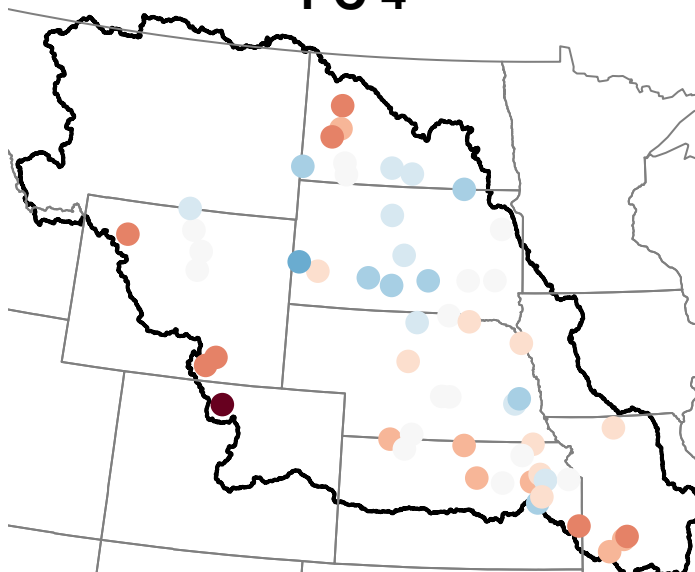
PC 2



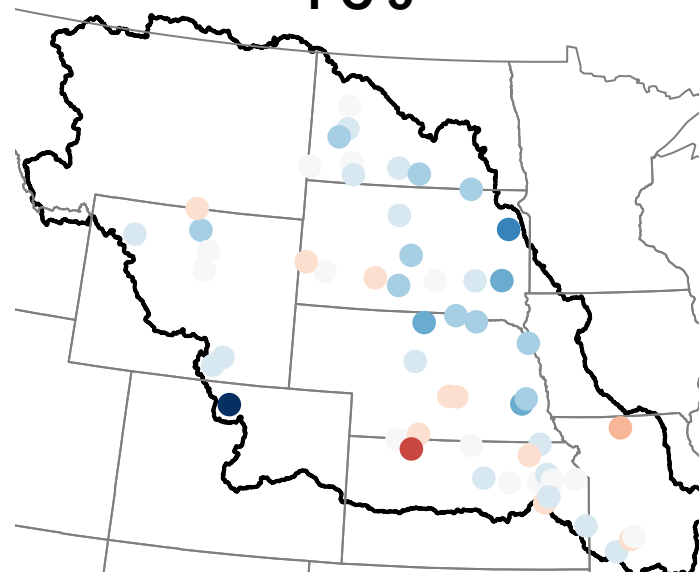
PC 3



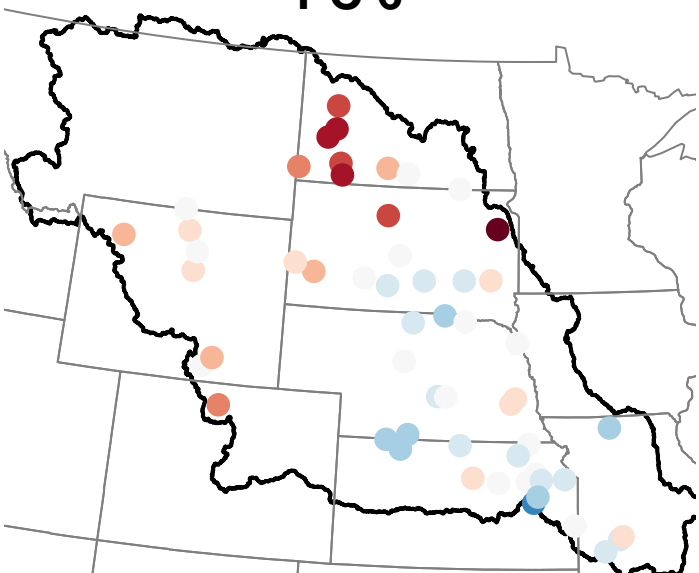
PC 4



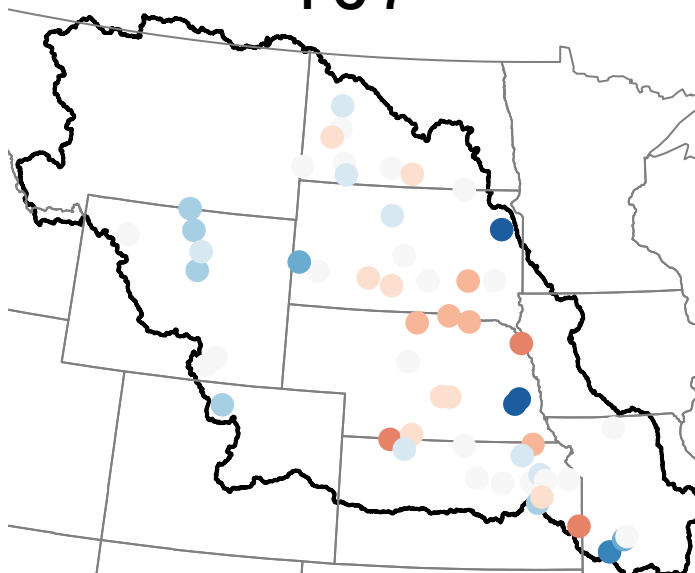
PC 5



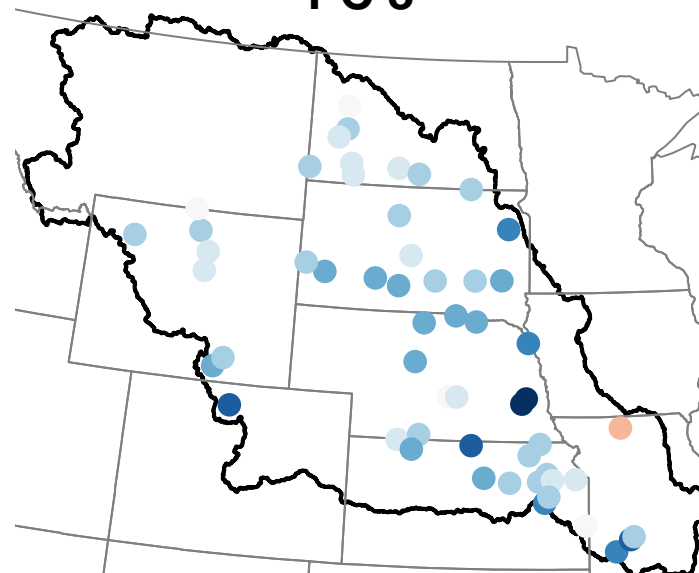
PC 6



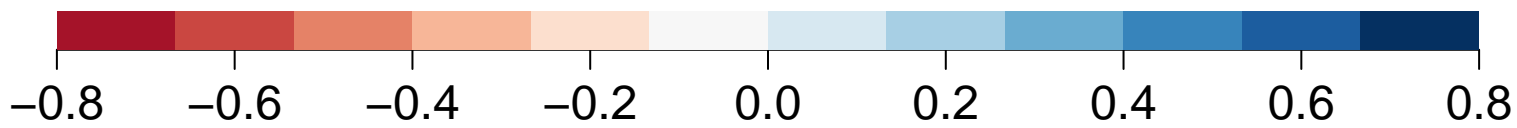
PC 7



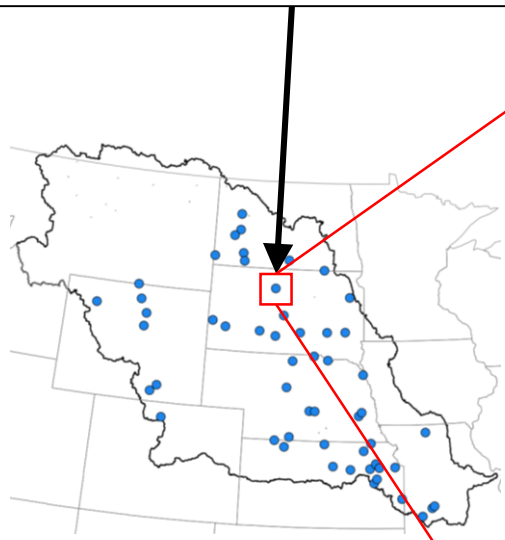
PC 8



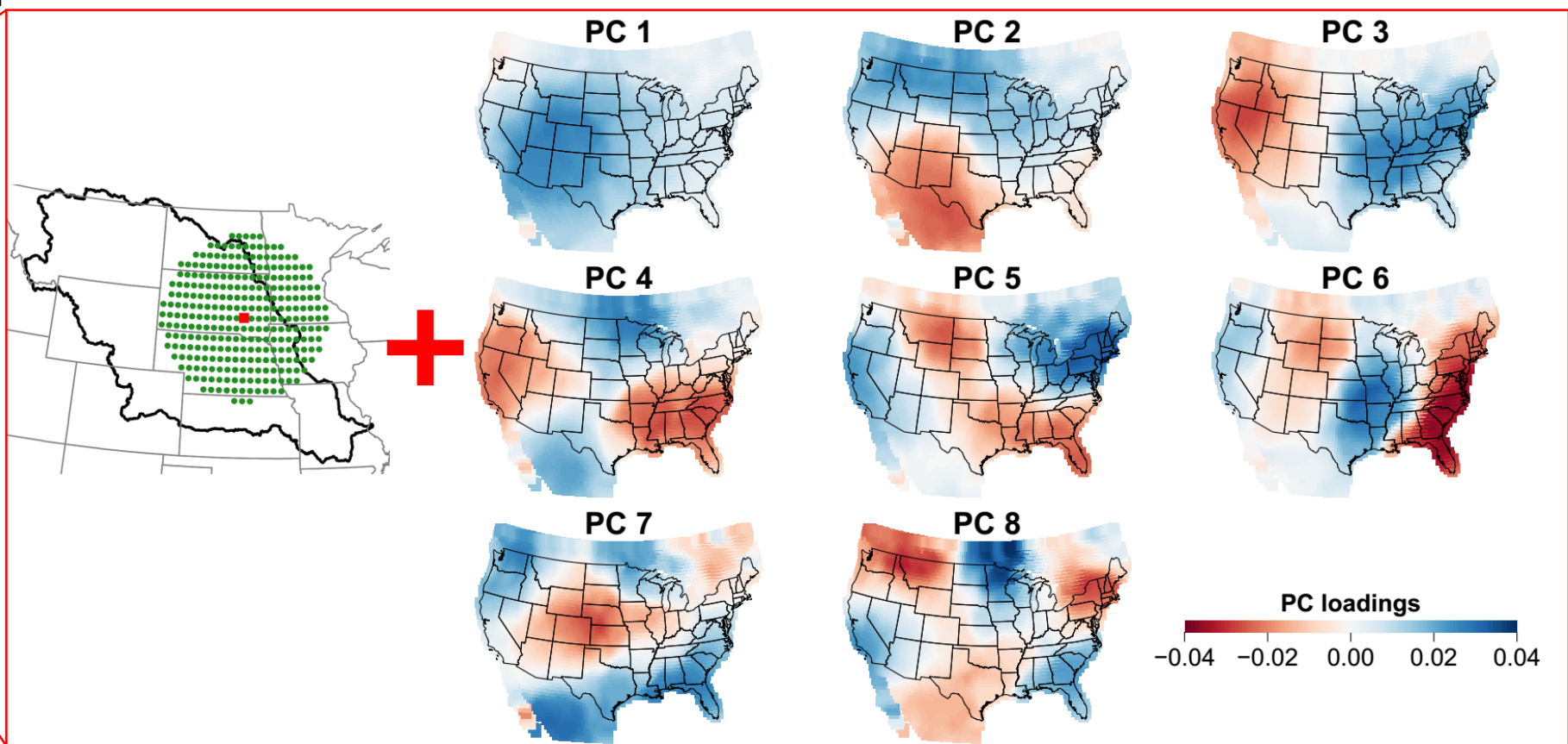
correlation coefficient



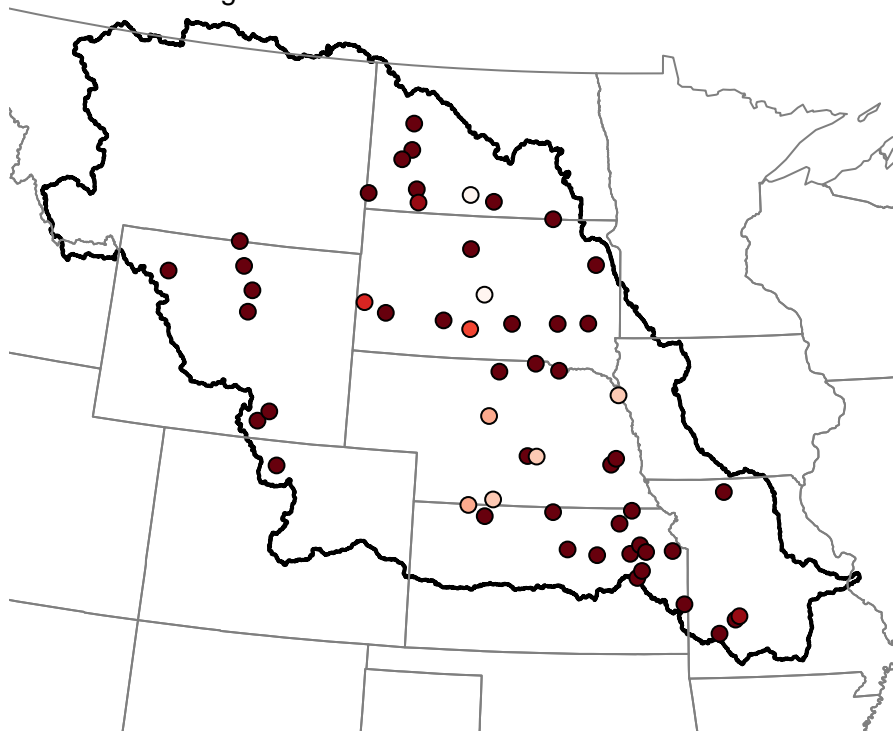
# Target streamflow



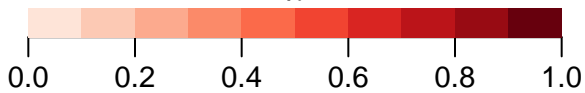
# Explanatory variables



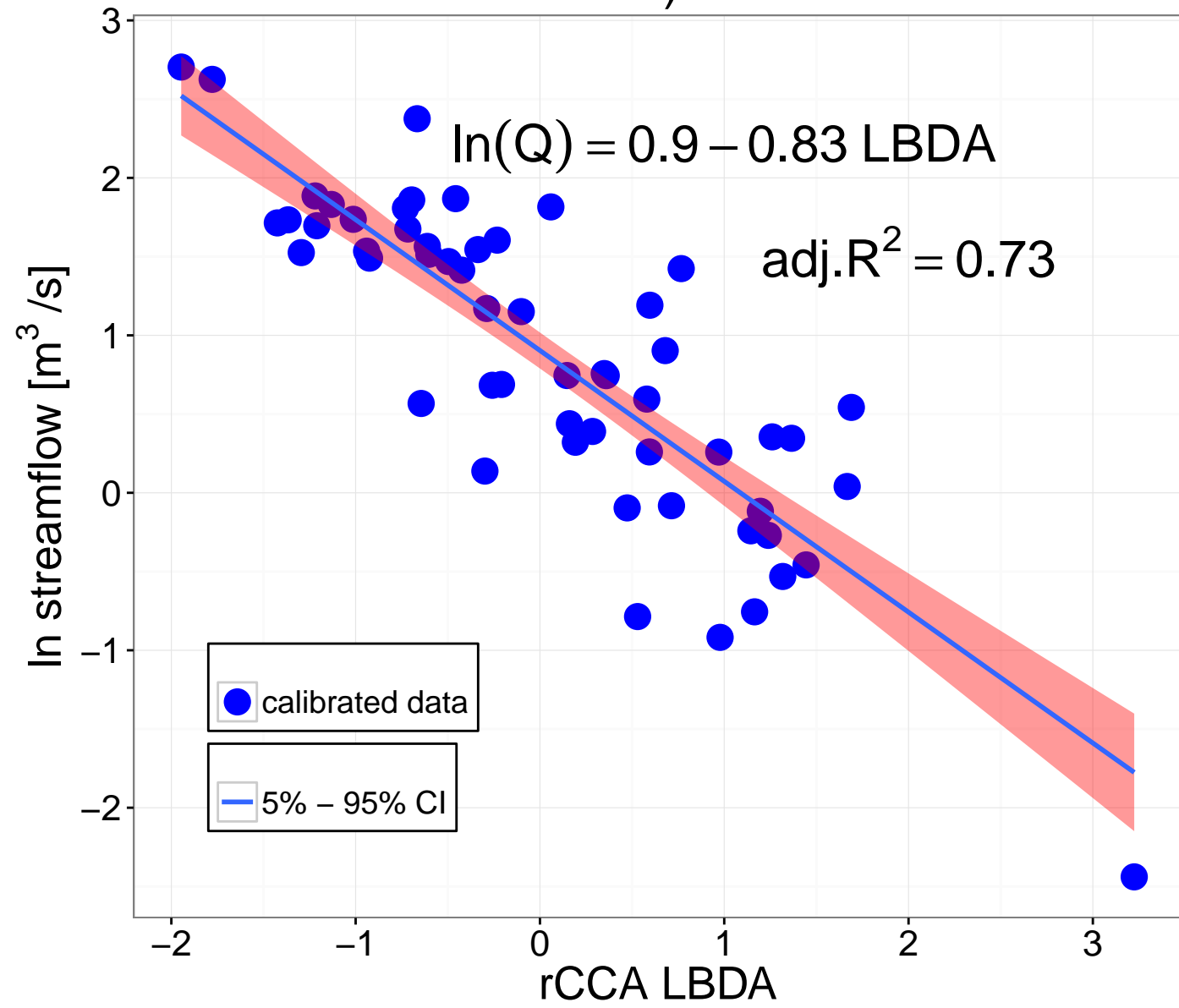
$\lambda$  for LBDA grids within a 450 km radius and US-wide PCs



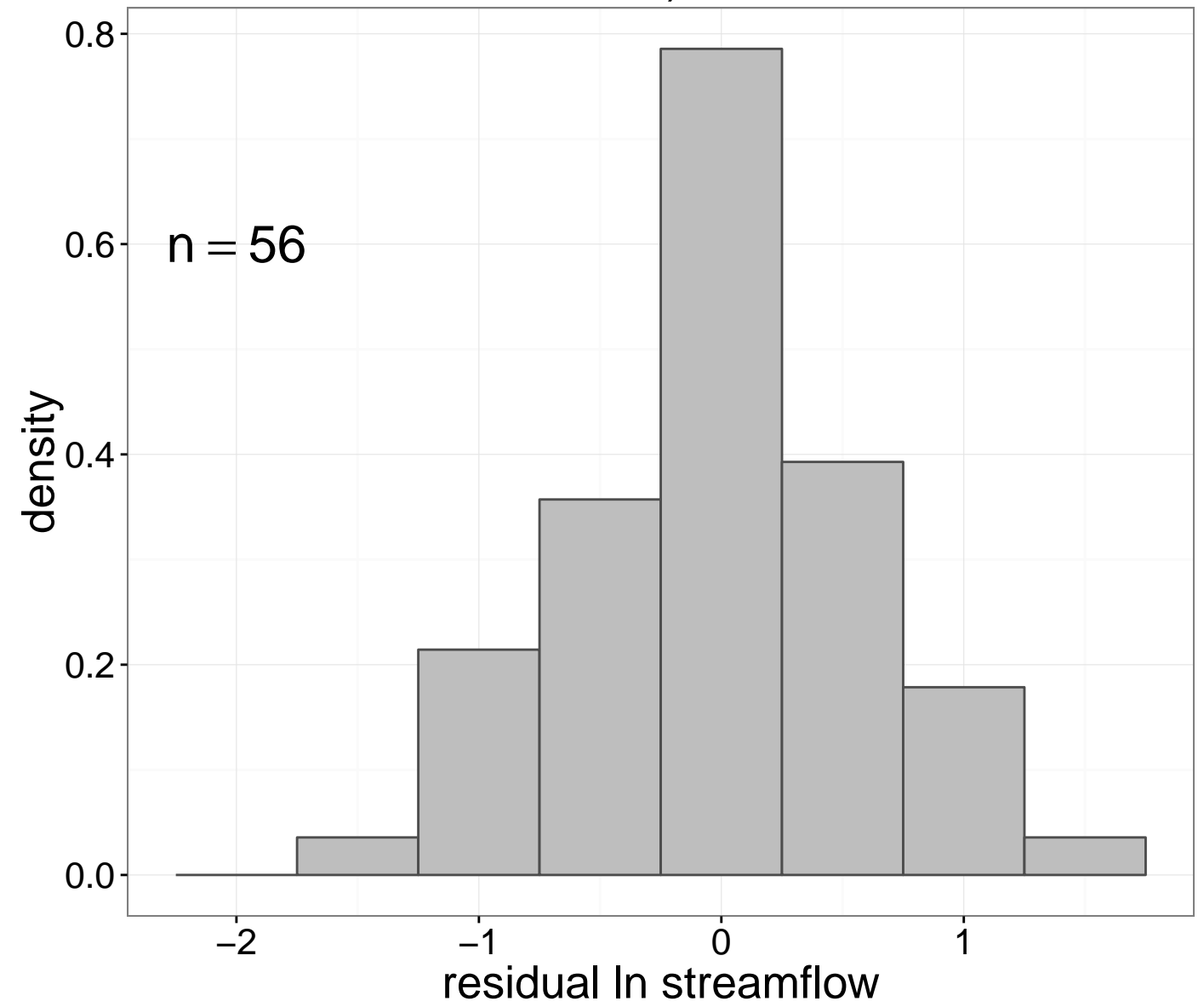
$\lambda$



a)

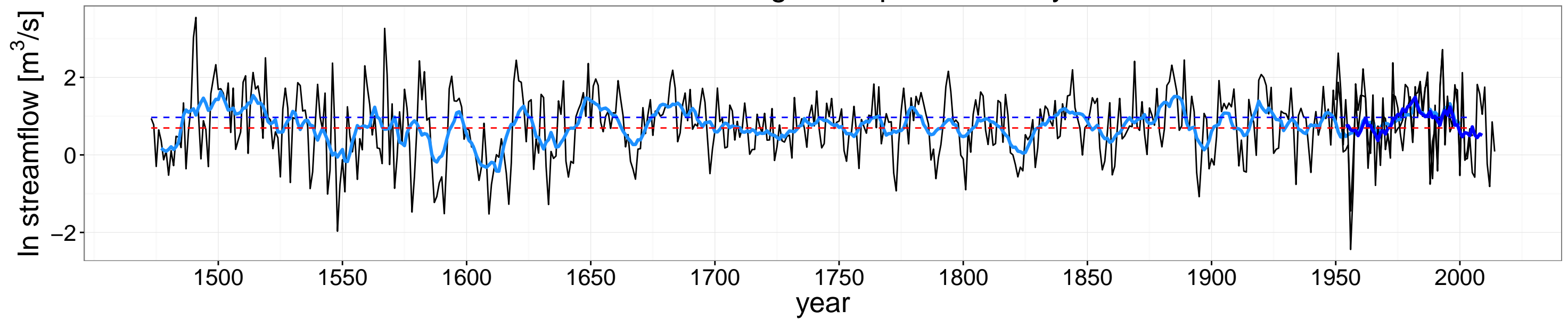


b)

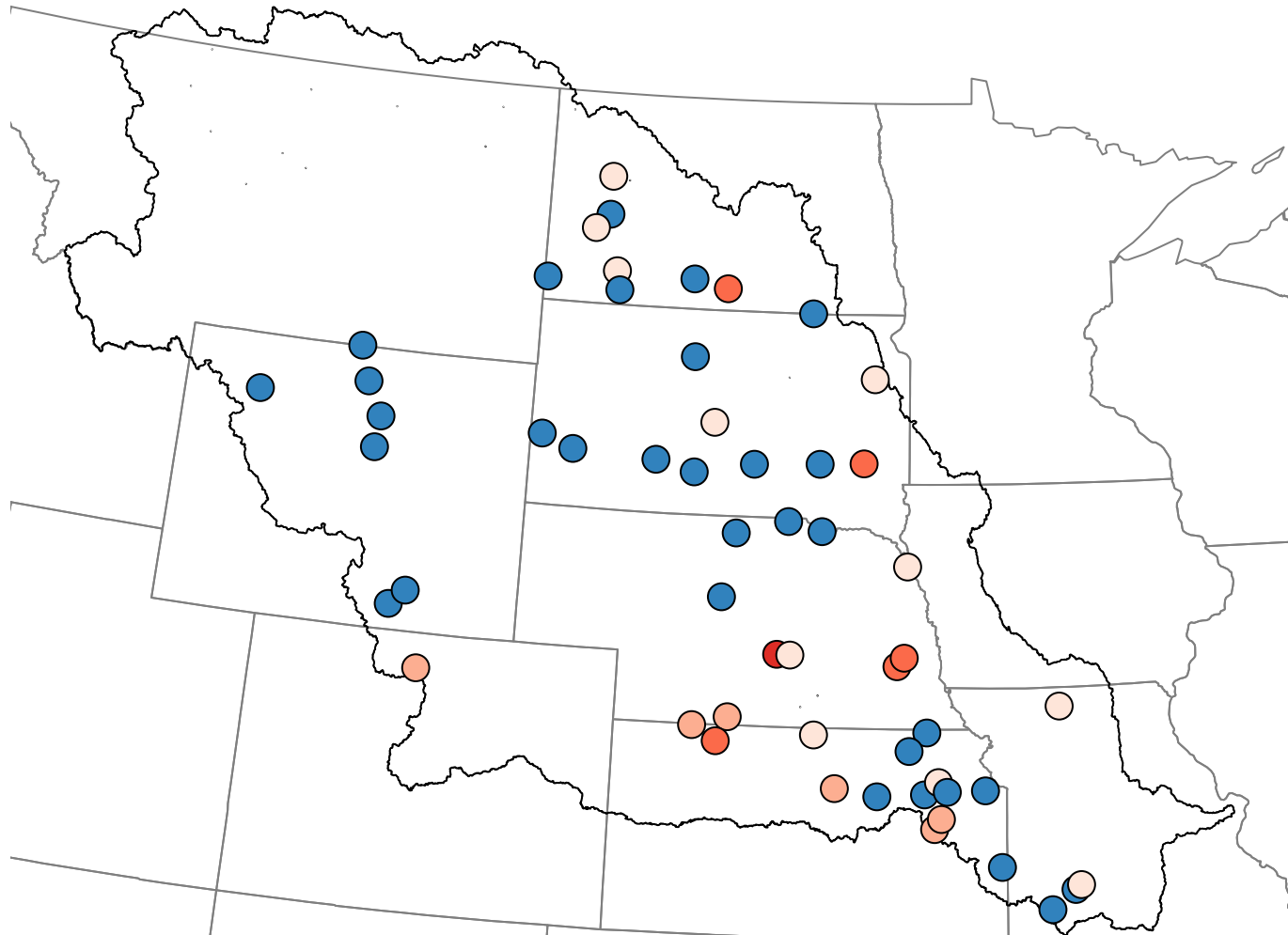


c)

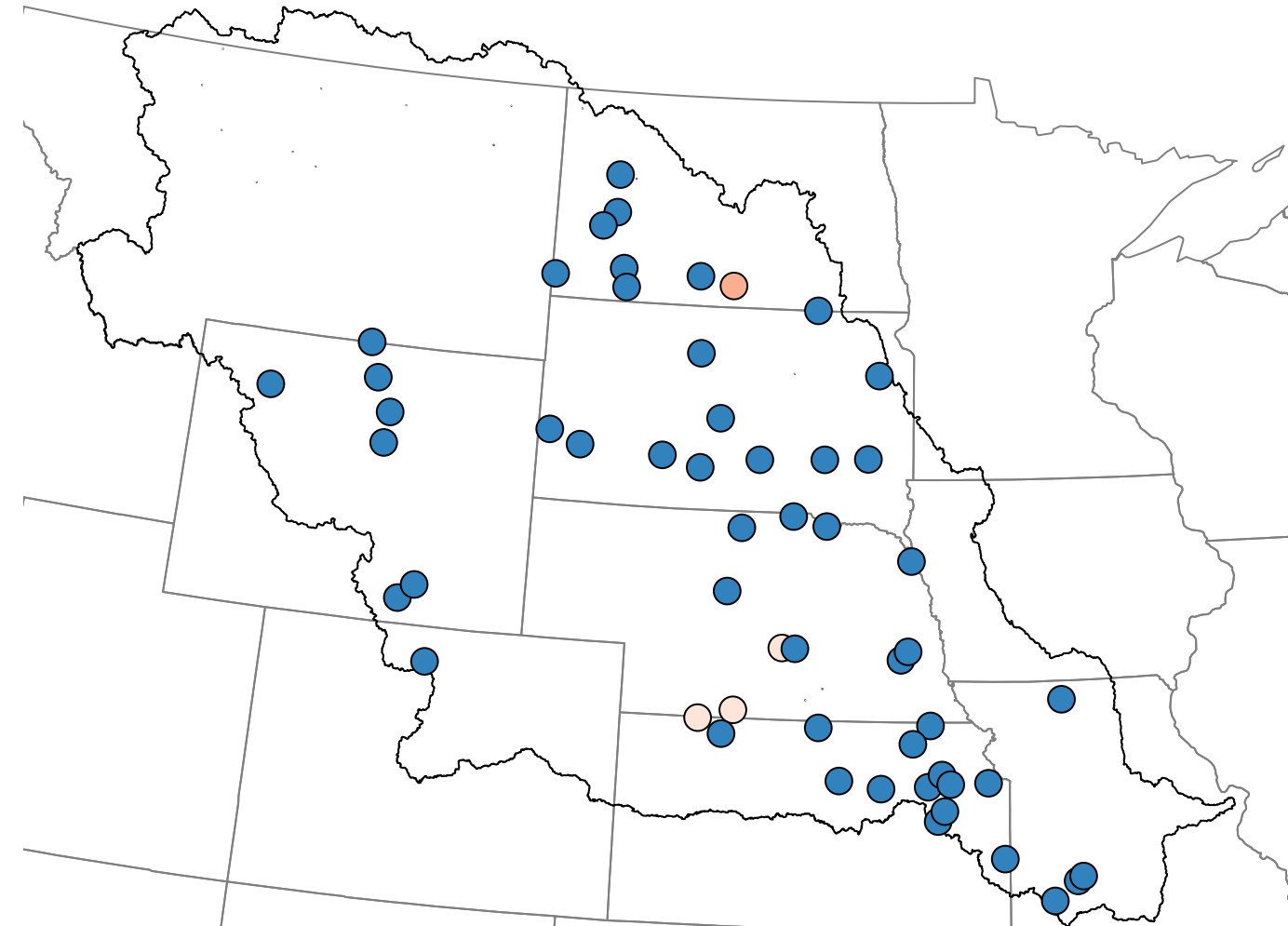
max duration flood inst/paleo: 14 / 34 years  
 max duration drought inst/paleo 7 / 23 years



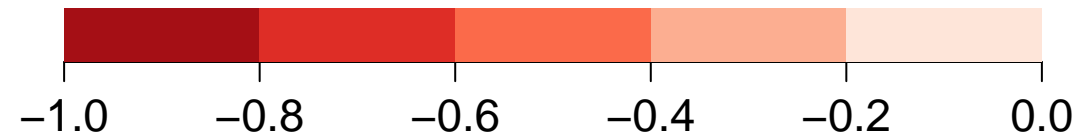
**a) Coefficient of efficiency (CE)**



**b) Reduction in error (RE)**

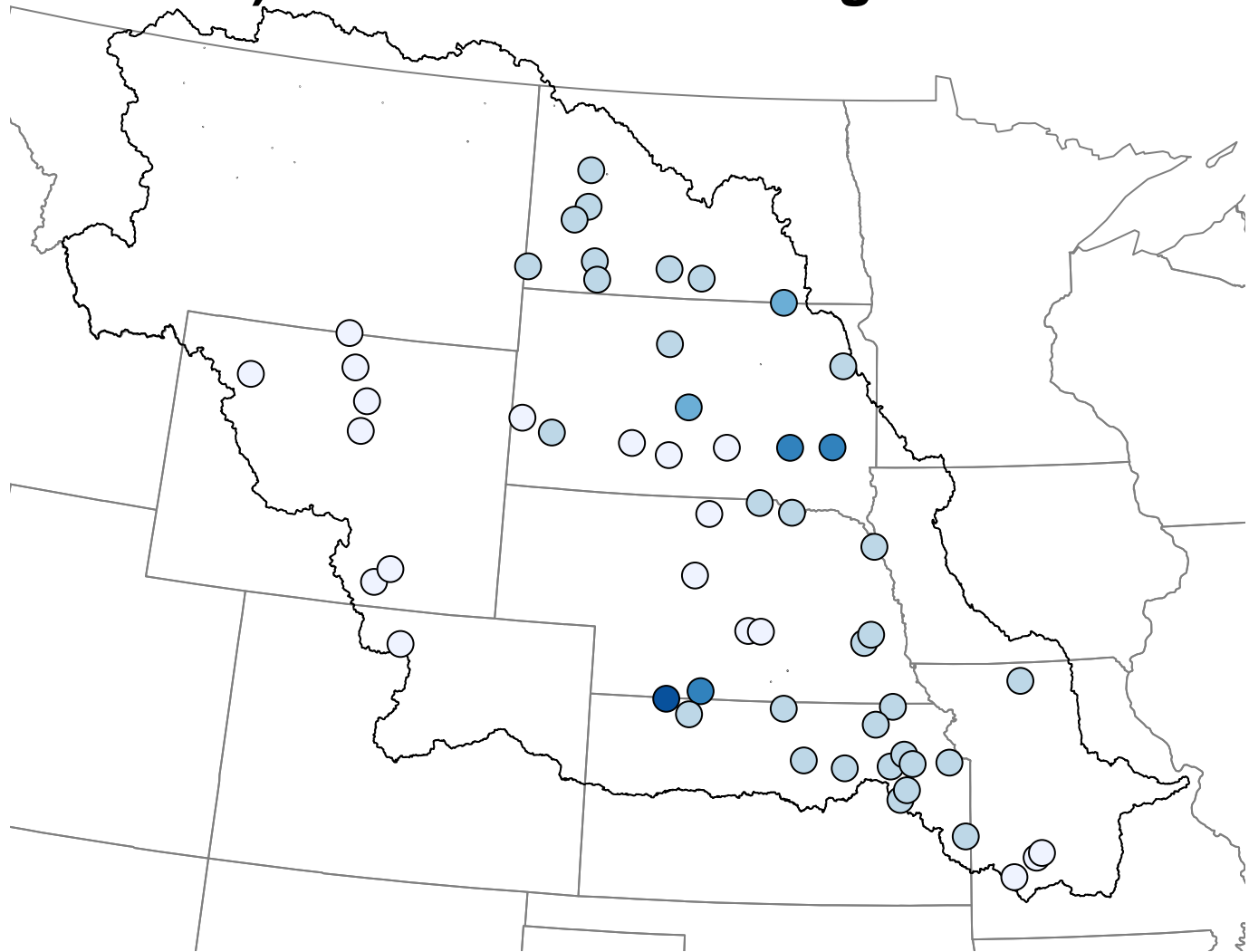


**CE and RE values**

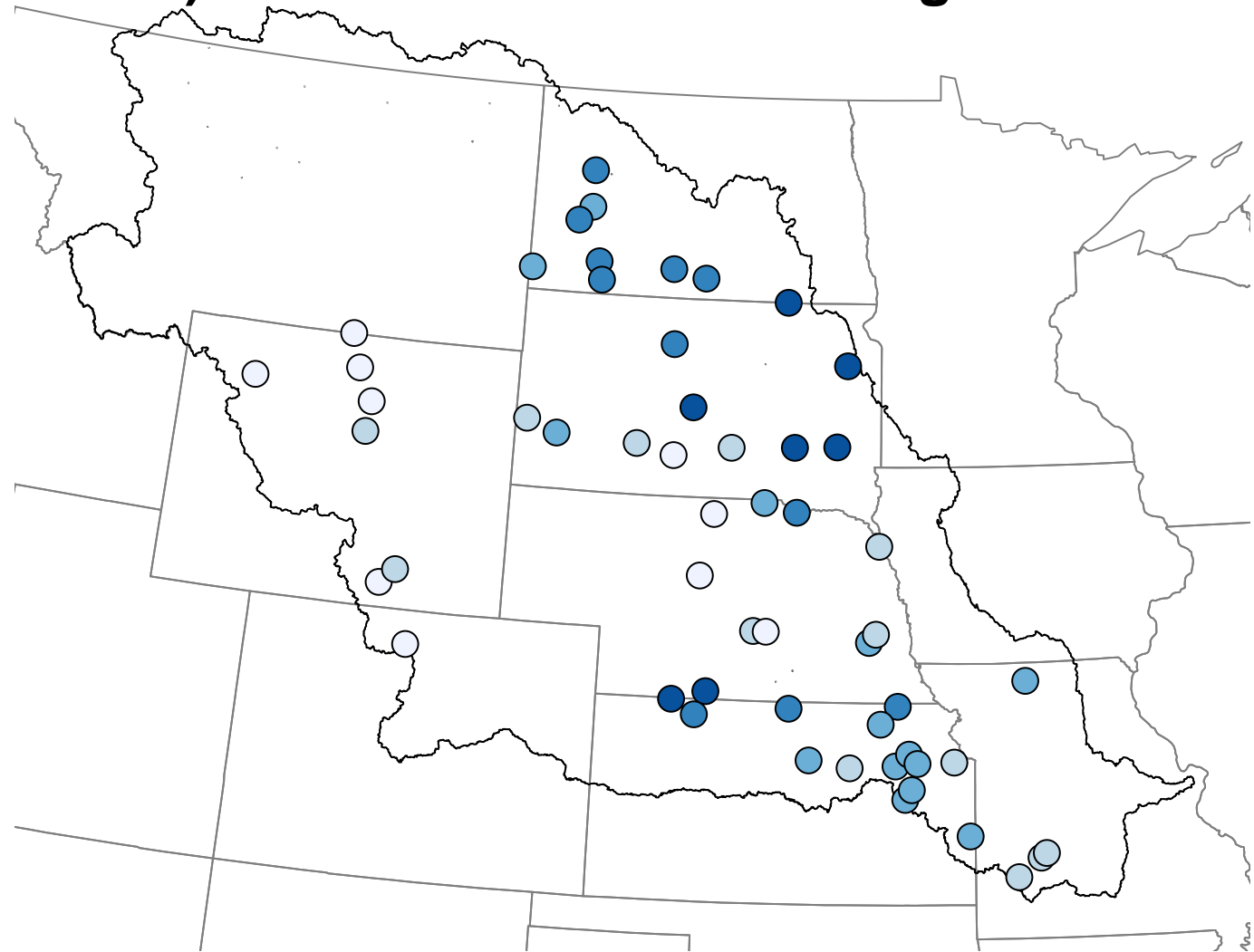


● positive CE

**a) calibrated natural log flows**

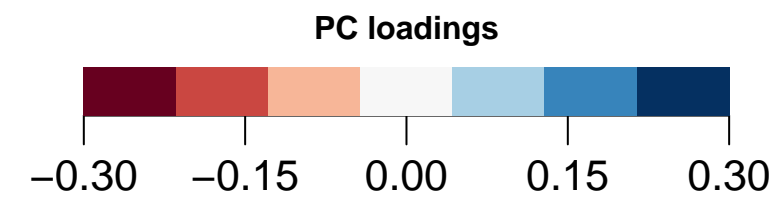
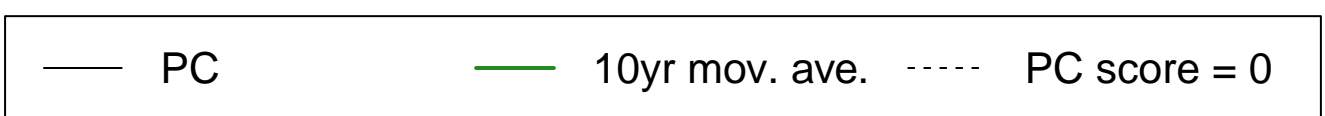
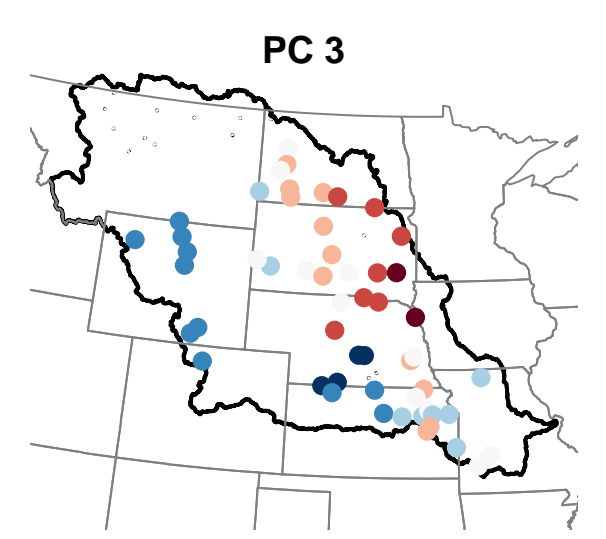
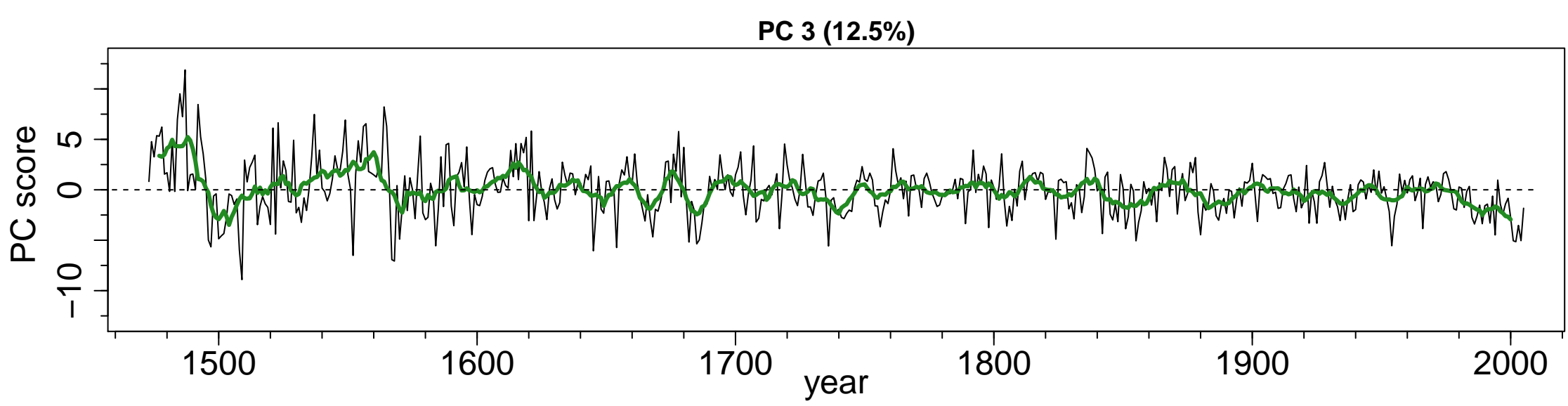
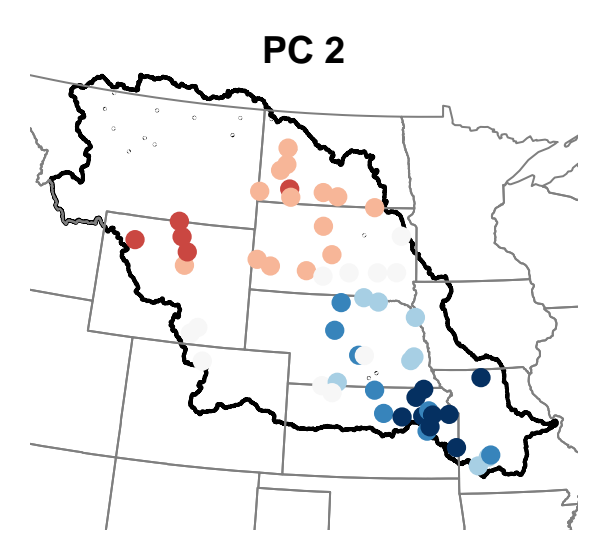
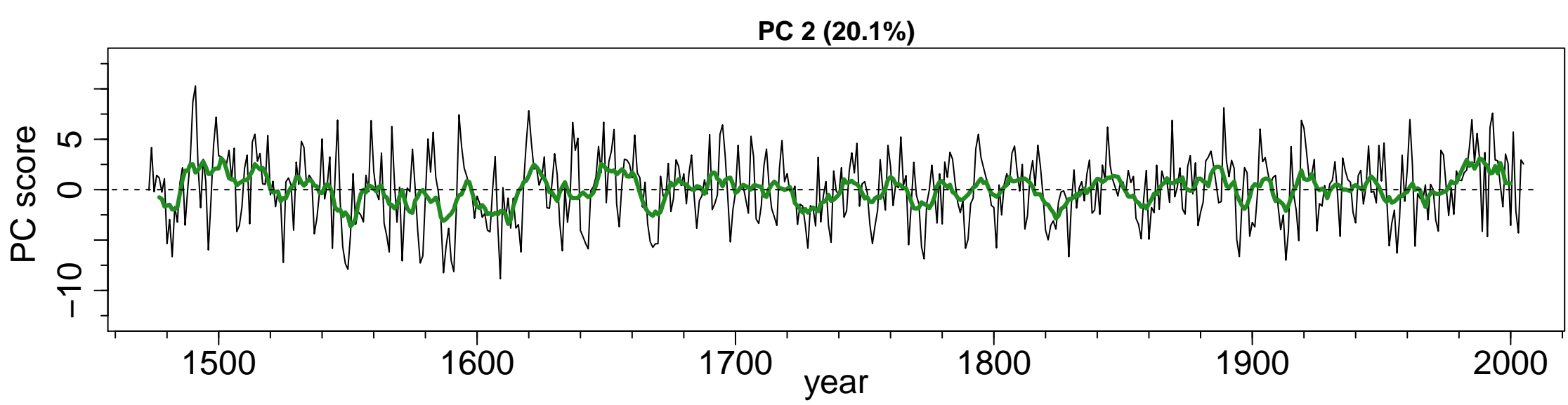
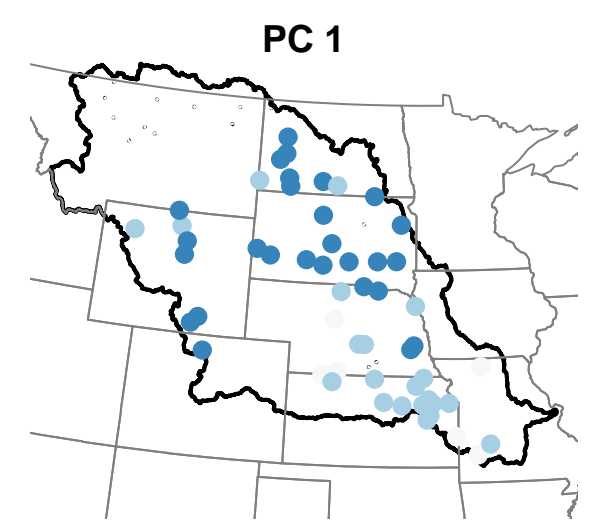
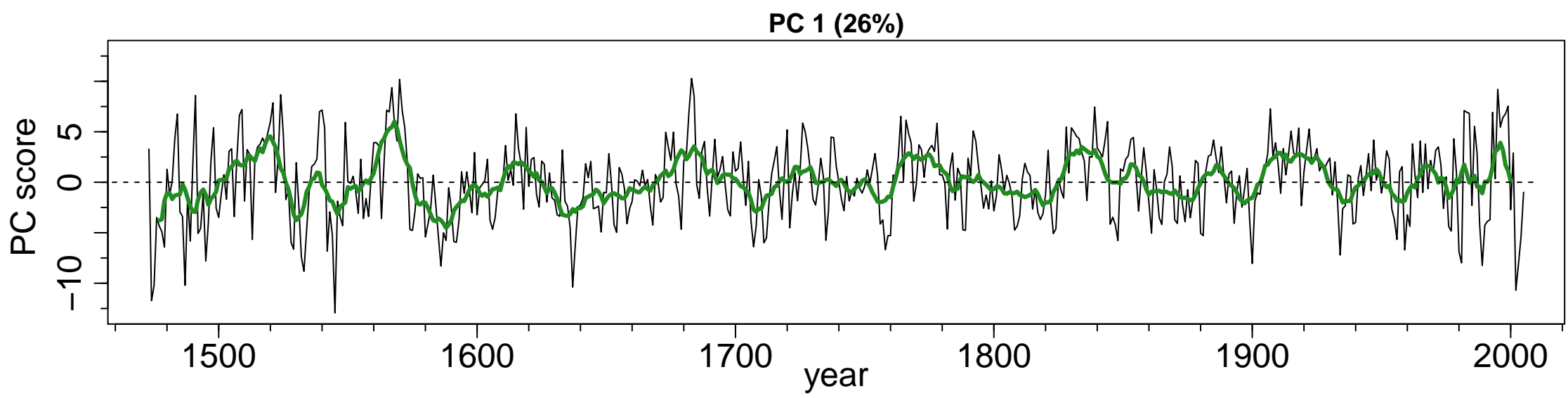


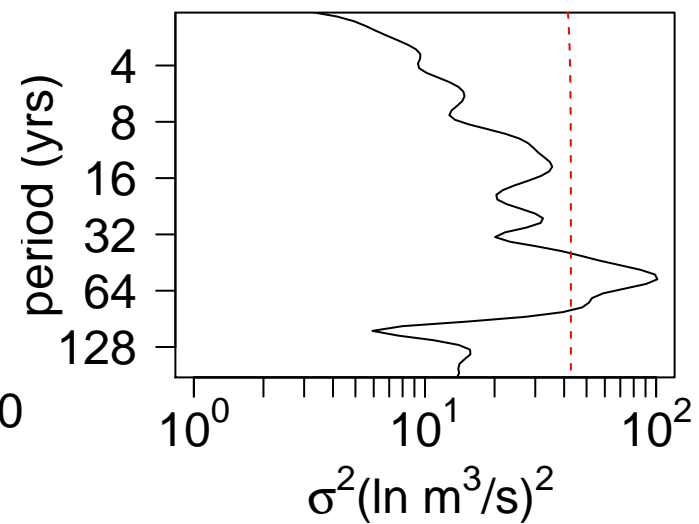
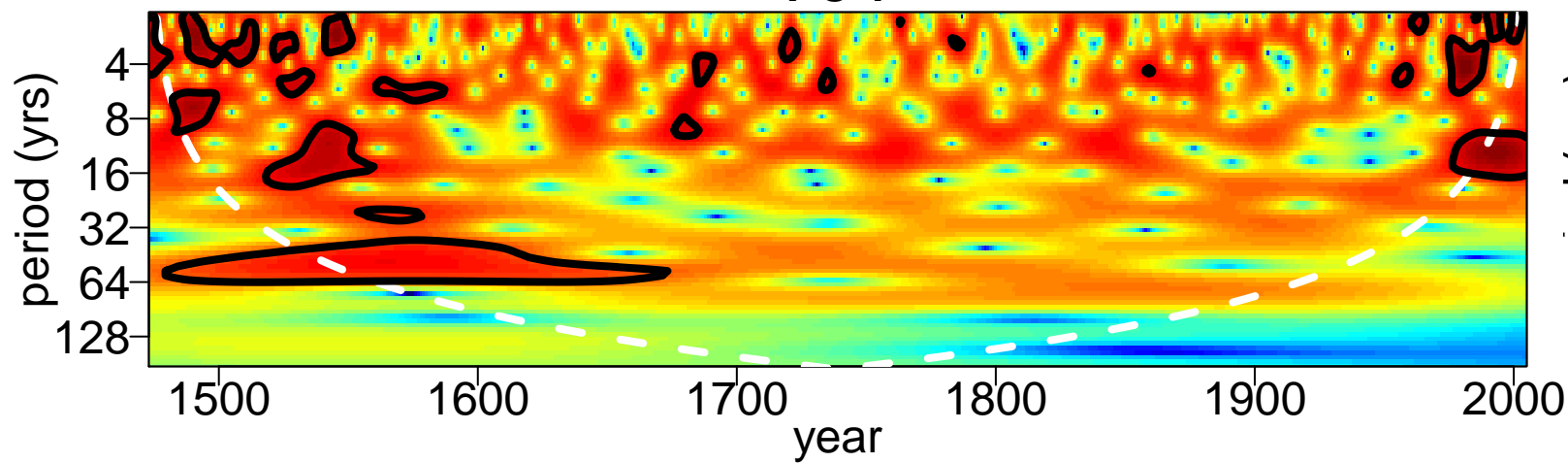
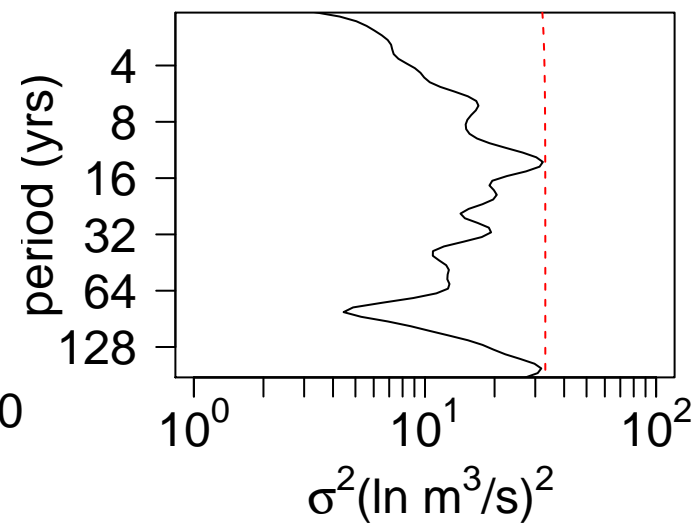
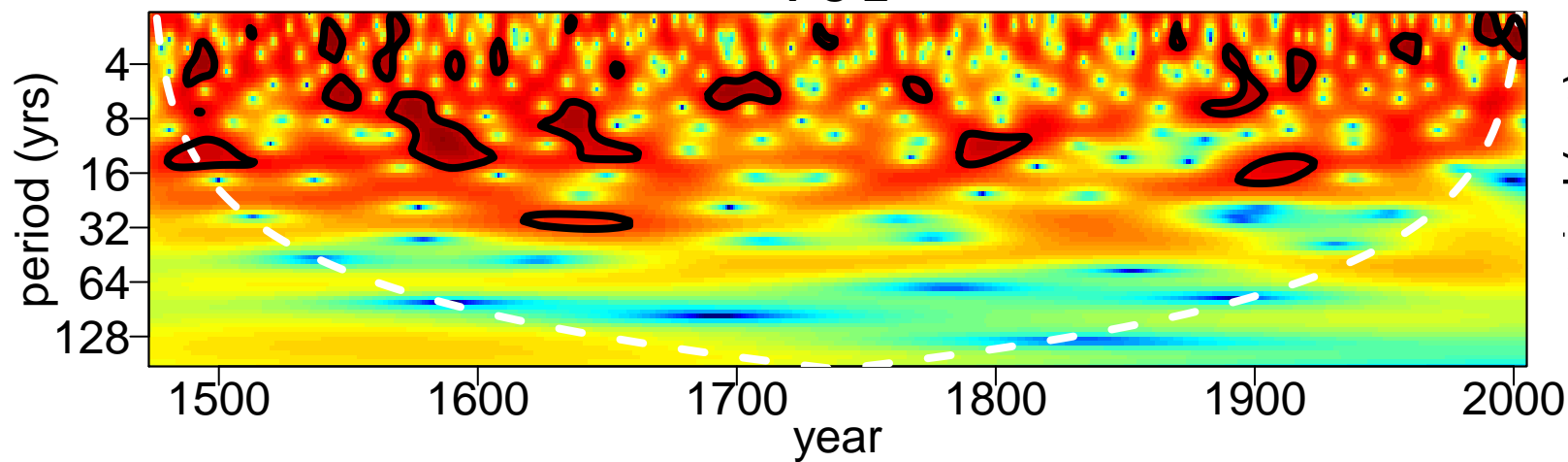
**b) cross validated natural log flows**



**median |residual|**





**PC 1****PC 2****PC 3**