A Wavelet-Based Approach to Streamflow Event Identification and Modeled Timing Error Evaluation

Erin Towler[1,*] and James L. McCreight[1,2]
[1] National Center for Atmospheric Research (NCAR), P.O. Box 3000, Boulder, CO 80307
[*] Corresponding author: towler@ucar.edu, https://orcid.org/0000-0002-1784-1346
[2] orcidid: 0000-0001-6018-425X

**Abstract**

Streamflow timing errors (in the units of time) are rarely explicitly evaluated, but are useful for model evaluation and development. Wavelet-based approaches have been shown to reliably quantify timing errors in streamflow simulations, but have not been applied in a systematic way that is suitable for model evaluation. This paper provides a step-by-step methodology that objectively identifies events, and then estimates timing errors for those events, in a way that can be applied to large-sample, high-resolution predictions. Step 1 applies the wavelet transform to the observations, and uses statistical significance to identify observed events. Step 2 utilizes the cross-wavelet transform to calculate the timing errors for the events identified in Step 1; this ~~. This step~~ ~~e approach also~~ includes the ~~a quantification of the~~ ~~confidence in~~diagnostic of model event "hits", and timing errors are only assessed for hits~~if the~~ ~~model "missed" observed evend and if the timing error estimates~~should not be considered. The methodology is illustrated using real and simulated stream discharge data from several locations to highlight key method features. The method groups event timing errors by dominant timescales, which can be used to identify the potential processes contributing to the timing errors and the associated model development needs.~~ ~~For instance, timing errors that are associated with the diurnal melt cycle are identified. The method is also useful for documenting and evaluating model performance in terms of defined standards. This is illustrated by showing version-over-version performance of the National Water Model (NWM) in terms of timing errors.

## 1. Introduction

Common verification metrics used to evaluate streamflow simulations are typically aggregated measures of model performance, e.g., the Nash Sutcliffe Efficiency (NSE) and the related root mean square error (RMSE). Although typically used to assess errors in amplitude, these statistical metrics include contributions from errors in both amplitude and timing (Ehret and Zehe 2011), making them difficult to use for diagnostic model evaluation (Gupta et al. 2008). Furthermore, common verification metrics are calculated using the entire time series, whereas timing errors require comparing localized features or events in the data. This paper focuses explicitly on event timing error estimation, which is not routinely evaluated, despite its potential benefit for model diagnostics (Gupta et al. 2008) and practical forecast guidance (Liu et al. 2011).

The fundamental challenge with evaluating timing errors is identifying what constitutes as an "event" in the two time series being compared. Identifying events is typically subjective, time consuming, and not practical for large-sample hydrological applications (Gupta et al. 2014). A variety of baseflow separation methods, ranging from physically-based to empirical, have been developed to identify hydrologic events (see Mei and Anagnostou 2015 for a summary), though many of these approaches require some manual inspection of the hydrographs. Merz et al. (2006) put forth an automated approach, but it requires a calibrated hydrologic model, which is a limitation in data poor regions. Koskelo et al. (2012) developed a simple, empirical approach that only requires rainfall and runoff time series, but is limited to small watersheds and daily data. Mei and Anagnostou (2015) introduce an automated physically-based approach, which is demonstrated for hourly data, though one caveat is that basin events need to have a clearly detectable recession period. ~~Most~~ Additional methods ~~for identifying events~~ have focused on

53 ~~flooding events. One common approach to~~ identifying flooding events ~~is to~~ us~~ing~~e peak-over-

54 threshold methods. The thresholds used for such analyses are often either based on historical

55 percentiles (e.g., the 95th percentile) or on local impact levels (river stage), such as the National

56 Weather Service (NWS) flood categories (NOAA National Weather Service, 2012). Timing error

57 metrics are often calculated from the peaks of these identified events. For example, the Peak

58 Time Error, or its derivative the Mean Absolute Peak Time Error, requires matching observed

59 and simulated event peaks, and calculating their offset (Ehret and Zehe 2011). While this may be

60 straightforward visually, it can be difficult to automate; some of the reasons for this are discussed

61 below.

62       Difficulties arise using thresholds for event identification. For example, exceedances can

63 cluster if a hydrograph vacillates above and below a threshold, begging the question: Is it one or

64 multiple events? Which peak should be used for the assessment? In the statistics of extremes,

65 declustering approaches can be applied to extract independent peaks (e.g., Coles 2001), but this

66 reductionist approach may miss relevant features. For instance, if background flows are elevated

67 for a longer period of time before and after the occurrence of these "events", the threshold-based

68 analysis identifies features of the flow separately from the primary hydrologic process

69 responsible for the event. If one focuses just on peak timing differences in this example, that

70 timing error may only apply to some small fraction of the total flow of the larger event which

71 happens mainly below the threshold. Further, for overall model diagnosis that focuses on model

72 performance for all events, not just flood events, variable thresholds would be needed to account

73 for different kinds of events (e.g., a daily melt event versus a convective precipitation event).

74       Using a threshold-approach to identify events and timing error assessment, Ehret and

75 Zehe (2011) develop an intuitive assessment of hydrograph similarity, the Series Distance. This

76    algorithm is later improved upon by Siebert et al. (2016). The procedure matches observed and

77    simulated segments (rise or recession) of an event, and then calculates the amplitude and timing

78    errors, as well as the frequency of event agreement. The Series Distance requires smoothing the

79    time series, identifying an event threshold, and selecting a time range to consider two segments

80    matching.

81          Liu et al (2011) developed a wavelet-based method for estimating model timing errors.

82    Although wavelets have been applied in many hydrologic applications such as model analysis

83    (e.g. Lane 2007; Weedon et al. 2015; Schaefli and Zehe 2009, Rathinasamy et al. 2014) and

84    post-processing (Bogner and Kalas 2007; Bogner and Pappenberger 2011), Liu et al. were the

85    first to use it for timing error estimation. Liu et al. (2011) apply a cross-wavelet transform

86    technique to streamflow time series for 11 headwater basins in Texas. Timing errors are

87    estimated for medium- to -high- flow "events" that are determined a priori by threshold

88    exceedance. They use synthetic as well as real streamflow simulations to test the utility of the

89    approach. They show that the technique can reliably estimate timing errors, though they

90    conclude that it is less reliable for multi-peak or consecutive "events" (defined qualitatively).

91    ElSaadani and Krajewski (2017) followed the cross-wavelet approach used by Liu et al (2011) to

92    provide similar analysis and further investigate the effect of the choice of mother wavelet on the

93    timing error analysis. Ultimately, they recommended that in the situation of multiple, adjoining

94    flow peaks the improved time localization of the Paul wavelet might justify its poorer frequency

95    localization compared the Morlet wavelet.

96          Liu et al. (2011) provide a starting point for the work in this paper where we develop two

97    new bases for their method: 1) objective event identification for timing error evaluation and 2)

98    the use of observed events as the basis for the model timing error calculations. The latter is

99    important for "model benchmarking", i.e., the practice of evaluating models in terms of defined

100   standards (e.g., Luo, et al. 2012; Newman et al. 2017). Here, the use of observed events provides

101   a baseline by which to evaluate changes and to compare multiple versions or experimental

102   designs.

103        This paper provides a methodology for using wavelet analysis to quantify timing errors in

104   hydrologic simulations. Our contribution is a systematic approach that integrates 1) statistical

105   significance to identify events with 2) a basis for timing error calculations independent of model

106   simulations (i.e., benchmarking). We apply our method to timing error evaluation of high-

107   resolution streamflow prediction. The paper is organized as follows: Section 2 describes the

108   observational and simulated data used. ~~provides an overview of the conceptual approach of using~~

109   ~~wavelets to identify events and estimate timing errors, and~~ Section 3 provides the detailed

110   methodology of using wavelets to identify events and estimate timing errors in a synthetic

111   example. ~~In Section 4, we describe the software and data, as well as provide a simple illustration~~

112   ~~of the method using real and simulated streamflow data.~~ In Section 4~~5~~, we ~~provide~~

113   ~~results~~demonstrate the method using real and simulated streamflow data for several use cases,

114   and then illustrate the application of the method for version-over-version comparisons. ~~,~~

115   ~~including select examples to highlight features of the method and version-over-version~~

116   ~~comparisons.~~ Section 6~~5~~ is the discussion and conclusions, including how specific

117   methodological choices may vary by application.

118   **2. Data**
119
120        The application of the methodology is illustrated using real and simulated stream discharge

121   (streamflow, m3/s) data ~~from~~at three U.S. Geological Survey (USGS) stream gage locations

122   ~~representing~~in three different geographic regions: Onion Creek at US Highway 183, Austin,

Texas, for the South Central region (Onion Creek, TX; USGS site number 08159000), Taylor River at Taylor Park, Colorado, for the Intermountain West (Taylor River, CO; USGS site number 09107000), and Pemigewasset River at Woodstock, New Hampshire, for New England (Pemigewasset River, NH; USGS site number 01075000). We use the USGS instantaneous observations averaged on an hourly basis.

NOAA's National Water Model (NWM, https://www.nco.ncep.noaa.gov/pmb/products/nwm/) is an operational model that produces hydrologic analyses and forecasts over the continental United States (CONUS) and Hawaii (as of version 2.0). The model is forced by downscaled atmospheric states and fluxes from NOAA's operational weather models. Next, the NoahMP (Niu et al 2011) land surface model calculates energy and water states and fluxes. Water fluxes propagate down the model chain through overland and subsurface (soil and aquifer representations) water routing schemes to reach a stream channel model. The NWM applies the three parameter Muskingum-Cunge river routing scheme to a modified version of the NHD-Plus version 2 (McKay et al. 2012) river network representation (Gochis et al 2020).

In this study, NWM simulations are taken from each version's retrospective runs (https://docs.opendata.aws/nwm-archive/readme.html). These are continuous simulations (not cycles) run for the period October 2010 to November 2016 and forced by the National Land Data Assimilation System (NLDAS)-2 product as atmospheric conditions. The nudging data assimilation was not applied in these runs either. We use NWM discharge simulations from versions V1.0, V1.1, and V1.2 (not all version may be publicly available).

144 ~~To apply the methodology, we note that the observed and simulated datasets must be paired~~

145 ~~(overlapping).~~ The methodology developed in this paper is implemented in the R language and is

146 made publicly available, as detailed in the code availability section at the end of the manuscript.

147 ~~**2. Conceptual Overview**~~

148 ~~This section provides the technical description of the methodology, and the steps can be seen~~

149 ~~in an accompanying flowchart (Supplemental Figure 1).~~~~Before going into technical details of the~~

150 ~~Method (Section 3), we provide a conceptual overview of the approach of using wavelets to~~

151 ~~identify events and estimate timing errors. We provide a nomenclature table (Supplemental Table~~

152 ~~1) of key terms relevant to the approach. The wavelet transform (WT) expands the dimensionality~~

153 ~~of the original time series by introducing the timescale (or period) dimension and returns power as~~

154 ~~a function of both time and timescale (e.g. Torrence and Compo, 1998). This is illustrated in Figure~~

155 ~~1: the streamflow time series (panel a) is expanded into a 2-dimensional wavelet power spectrum~~

156 ~~(panel b). Where traditional model errors, such as the aforementioned RMSE or NSE, reduce the~~

157 ~~information of the time series to a single statistic, wavelet analysis expands the input signal and~~

158 ~~provides information on the dominant timescales of the time series at each time. Wavelet analysis~~

159 ~~can therefore detect localized signals in time series (Daubechies 1990), including hydrologic time~~

160 ~~series, which are often irregular or aperiodic (i.e., events may be isolated and don't regularly~~

161 ~~repeat) or non-stationary. We note that in many wavelet applications, timescale is referred to as~~

162 ~~"period". To emphasize that our study is more focused on irregular events and less on periodic~~

163 ~~behavior of time series, we use the term "timescale". The wavelet transform is the foundation of~~

164 ~~the view in this paper that events have characteristics of both time and timescale. Timing errors,~~

165 ~~calculated from events defined this way, therefore have dimensions of both time and timescale as~~

166 ~~well.~~

167    In their seminal wavelet study, Torrence and Compo (1998) outline a method for objectively

168    identifying statistical significance in the wavelet transform. We adopt this approach and define

169    "events" in the observed time series via statistical significance of the wavelet power spectrum. The

170    details are provided in the next section, however Figure 1 illustrates that the events in the input

171    time series (panel a) are defined as regions of the wavelet power spectrum shown in panel b: events

172    are inside the black contours (>= 95% confidence level) but not inside the cone of influence

173    (regions where the colors are muted, this is explained in detail in Section 3). The wavelet power

174    spectrum is only shown for the events in panel c. Events defined in this way are a function of both

175    time and timescale. Note that at a given time, events of different timescales can occur

176    simultaneously. What one may subjectively interpret as a single event in the input time series is

177    generally quantified by this definition as multiple coincident events at a variety of timescales each

178    with a different power (e.g. Figure 1, panel c). Although for some locations there may be physical

179    reasons to expect certain timescales to be important (e.g., seasonal cycle of snowmelt), the most

180    important scales at which hydrologic signals occur at a particular location are not necessarily

181    known a priori. The wavelet power can be examined across events to identify the most dominant,

182    or what we call "characteristic" timescales for a given time series; the procedure for this is detailed

183    later in the technical methodological section (Section 3.1.3). This approach to event detection is

184    objective, data-driven, and portable across diverse locations, which is important for large-sample

185    hydrologic applications. We point out that in the objective identification of events, we are not

186    limited to flooding events. Rather, events are defined more broadly: an event is when the wavelet

187    power falls outside its standard statistical power. This can be further subset into flooding events if

188    desired.

189    Once observed events are identified by the method, we can calculate timing errors

190    between observed and simulated time series.  The cross-wavelet timing error approach of Liu et

191    al (2011) is used, but we restrict our calculation of timing errors to the aforementioned regions of

192    statistically significant wavelet power in the observations; i.e., we calculate timing errors in

193    terms of *observed* events (Figure 1c). Because both the phase (timing error) and the significance

194    of the cross wavelet transform (XWT) computed between the observed and modeled time series

195    depends on the modeled time series, we use the observed event definition (Figure 1c) in the

196    calculation of the timing errors to provide a common, consistent basis independent of the models

197    evaluated (i.e., benchmarking). The portions of the observed wavelet spectrum used for

198    comparison may further be restricted depending on the analysis goals.

199    **3. Method**ology ~~for evaluating event timing errors~~

200

201    This section provides the technical description of the methodology, and the steps

202    can be seen in an accompanying flowchart (Supplemental Figure 1). This section provides the

203    description of the methodology using wavelets to identify events and estimate timing errors. The

204    steps can be seen in an accompanying flowchart (Figure 1) and nomenclature table (Table 1),

205    which defines key terms of the approach. To facilitate understanding, the steps are illustrated

206    accompanied by an application of the methodology to an observed time series of an isolated peak

207    in Onion Creek, TX (Figure 2a), (figure 2a) and the synthetic modeled time series which is

208    identical to the observation time seriesuniformly but shifted 5 hours in to the future (figure 3a, note

209    the log scalescale).

210

211    *3.1. Step 1. Identify observed events*

212

213    The first step towards evaluating timing errors is to identify a set of observed events for

214    which the timing error should be calculated. We break this step into three sub-steps: 1a. Apply the

215    wavelet transform to observations, 1b. Determine all observed events using significance testing,

216    and 1c. Sample observed events to an event-set relevant to analysis. ~~To facilitate understanding,~~

217    ~~the steps are accompanied by an application of the methodology to an observed time series of an~~

218    ~~isolated peak in Onion Creek, TX (Figure 2a).~~

219    3.1.1. Step 1a. Apply wavelet transform to observations

220    First, we apply the continuous wavelet transform (WT) to the observed time series. The

221    main steps and equations for the WT are provided here, though the reader is referred to Torrence

222    and Compo (1998) and Liu et al. (2011) for more details.

223    Before applying the WT, a mother wavelet needs to be selected. In Torrence and Compo

224    (1998), they discuss the key factors that should be considered when choosing the mother

225    wavelet. There are four main considerations, including (i) orthogonal or nonorthogonal, (ii)

226    complex or real, (iii) width, and (iv) shape. In this study, we follow Liu et al. (2011) in selecting

227    the nonorthogonal and complex Morlet wavelet:

228    $$\psi(n) = \pi^{-1/4} e^{iw_0 n} e^{-n^2/2},$$
229

230    where $w_0$ is the non-dimensional frequency, with a value of 6 (Torrence and Compo, 1998).

231    Once the mother wavelet is selected, the WT is applied to a time series $x_n$, where n goes

232    from n=0 to n=N-1, with a time step of $\delta t$. The WT is the convolution of the time series with the

233    mother wavelet that has been scaled and normalized:

234    $$W_n(s) = \sum_{n'=0}^{N-1} x_{n'} \psi^* \left[ \frac{(n'-n)\delta t}{s} \right],$$
235

236    where $n'$ is the localized time in [0, N-1], $s$ is the scale parameter, and the asterix indicates the

237    complex conjugate of the wavelet function. The wavelet power is defined as $|W_n^2|$, which

10

238  represents the squared amplitude of an imaginary number when a complex wavelet is used as in

239  this study. We use the bias corrected wavelet power spectrum (Liu et al. 2007; Veleda et al.

240  2012), which ensures spectral peaks are power is comparable across timescales. We also identify

241  a maximum timescale *a priori* that corresponds to our application. We select 256 hours (~10

242  days), but this number could be higher or lower for other applications and there are no real

243  penalties for using too high a maximum (lower than the annual cycle).

244  The wavelet transform (WT) expands the dimensionality of the original time series by

245  introducing the timescale (or period) dimension. Wavelet power and returns power as is also a

246  function of both time and timescale (e.g. Torrence and Compo, 1998). This is illustrated in

247  Figure 12: the streamflow time series (panel a) is expanded into a 2-dimensional (2-D) wavelet

248  power spectrum (panel b). Where traditional model errors, such as the aforementioned RMSE or

249  NSE, reduce the information of the time series to a single statistic, wavelet analysis expands the

250  input signal and provides information on the dominant timescales of the time series at each time.

251  Wavelet analysis can therefore detect localized signals in time series (Daubechies 1990),

252  including hydrologic time series, which are often irregular or aperiodic (i.e., events may be

253  isolated and don't regularly repeat) or non-stationary. We note that in many wavelet applications,

254  timescale is referred to as "period". To emphasize that our study is more focused on irregular

255  events and less on periodic behavior of time series, we use the term "timescale". The wavelet

256  transform is the foundation of the view in this paper that events have characteristics of both time

257  and timescale. Timing errors, calculated from events defined this way, therefore have dimensions

258  of both time and timescale as well. We note that in many wavelet applications, timescale is

259  referred to as "period" and this axis is indeed the Fourier period in our plots. However, to

11

260 emphasize that our study is more focused on irregular events and less on periodic behavior of

261 time series, we use the term "timescale" to denote Fourier period (and not wavelet scale).

262 ~~We provide an overview of the main steps and equations for the wavelet transform here, though~~

263 ~~the reader is referred to Torrence and Compo (1998) and Liu et al. (2011) for more details.~~

264

265 ~~Before applying the WT, a mother wavelet needs to be selected. In Torrence and Compo~~

266 ~~(1998), they discuss the key factors that should be considered when choosing the mother~~

267 ~~wavelet. There are four main considerations, including (i) orthogonal or nonorthogonal, (ii)~~

268 ~~complex or real, (iii) width, and (iv) shape. In this study, we follow Liu et al. (2011) in selecting~~

269 ~~the nonorthogonal and complex Morlet wavelet:~~

270 $$\psi(n) = \pi^{-1/4} e^{iw_0 n} e^{-n^2/2},$$
271 ~~where $w_0$ is the non-dimensional frequency, with a value of 6 (Torrence and Compo, 1998).~~

272 ~~Once the mother wavelet is selected, the WT is applied to a time series $x_n$, where n goes~~

273 ~~from n=0 to n=N-1, with a time step of $\delta t$. The WT is the convolution of the time series with the~~

274 ~~mother wavelet that has been scaled and normalized:~~

275 $$W_n(s) = \sum_{n'=0}^{N-1} x_{n'} \psi^* \left[ \frac{(n'-n)\delta t}{s} \right],$$
276 where s is the scale parameter, the asterix indicates the complex conjugate of the wavelet

277 function. The wavelet power is defined as $|W_n^2|$. We use the bias corrected wavelet power

278 spectrum (Liu et al. 2007; Veleda et al. 2012), which ensures spectral peaks are comparable

279 across timescales. We also identify a maximum timescale that corresponds to our application.

280 We select 256 hours (~10 days), but this number could be higher or lower for other applications

281 and there are no real penalties for using too high a maximum (lower than the annual cycle).

282 ~~Because we are applying the WT to a finite time series, there are timescale-dependent~~ Because

283 we are applying the WT to a finite time series, there are timescale-dependent errors at the

284 beginning and end times of the power spectrum, where the entirety of the wavelet at each scale is

285 not fully contained within the time series. This region of the WT is referred to as the cone of

286 influence or COI (Torrence and Compo, 1998). ~~errors at the beginning and end times of the~~

287 ~~power spectrum. These are referred to as the cone of influence or COI (Torrence and Compo,~~

288 ~~1998).~~ Figure 2b illustrates the COI as the regions where the colors are muted; ~~W~~we ignore all

289 results within the COI in this study. ~~The details are provided in the next section, however Figure~~

290 ~~1 illustrates that the events in the input time series (panel a) are defined as regions of the wavelet~~

291 ~~power spectrum shown in panel b: events are inside the black contours (>= 95% confidence~~

292 ~~level) but not inside the cone of influence (regions where the colors are muted, this is explained~~

293 ~~in detail in Section 3).~~

294 We make several additional notes on the wavelet power and its representation in the

295 figures. The units of the wavelet power are those of the time series variance (m6/s2 for

296 streamflow) and it is natural to want to cast the power in a physical light or relate it to the time

297 series variance. Indeed, the power is often normalized by the time series variance when presented

298 graphically. However, it must be noted that the wavelet convolved with the time -series frames

299 the resulting power in terms of itself at a given scale. Wavelet power is a (normalized) measure

300 of how well the wavelet and the time series match at a given time and scale. The power can only

301 be compared to other values of power resulting from a similarly constructed WT. There are

302 various transforms that can be applied to aid graphical interpretation of the power (log, variance

303 scaling), but the utility of these often depends on the nature of the individual time series

304 analyzed. For simplicity, we plot the raw bias-rectified wavelet power in this paper.

305

306 3.1.2. Step 1b. Determine all observed events using significant testing

~~Once the WT is applied, the 2-dimensional (2-D) wavelet power spectra shows how the~~

~~features of the time series vary with both time and timescale.~~ In their seminal wavelet study,

Torrence and Compo (1998) outline a method for objectively identifying statistical significance

in the wavelet ~~transform.~~power ~~We adopt this approach and define "events" in the observed time~~

~~series via statistical significance of the wavelet power spectrum.~~ ~~To identify areas of~~

~~significance, we apply Torrence and Compo's (1998) approach that~~by compar~~es~~ing the

wavelet~~WT~~ power spectra with a power spectra from a red noise process. Specifically, the

observed time series is fitted with an order 1 autoregressive (AR1, or red noise) model, and the

WT is applied to the AR1 time series. The power spectr~~um~~a of the AR1 model provide~~s~~ the basis

for the statistical significance testing. Significance is determined if the power spectra are

statistically different using a chi-squared test~~ with 95% confidence~~.

~~We apply this here, and~~ Figure 2b shows significant (>= 95% confidence level) regions

of wavelet power~~illustrates the events, which are~~ inside ~~the~~ black contours ~~(>= 95% confidence~~

~~level) but not inside the COI.~~ ~~S~~Statistical significance indicates ~~an "event" at a given time and~~

~~timescale: that is, the~~ wavelet ~~that~~ power that falls outside the time series~~its~~ ~~standard~~ background

statistical power based on an AR1 model of the time series. Statistical significance of the wavelet

power can be thought of as events in the wavelet domain. We define events as regions of

significant wavelet power outside the COI. Figure 2c displays the wavelet power for the events

in this this time series~~The result is the set of all events, i.e., each event is a combination of time~~

~~and timescale (i.e., locations on the 2-D grid).~~ ~~The wavelet power spectrum is only shown for~~

~~the events in Figure 2~~panel c. ~~We emphasize that e~~Events defined in this way are a function of

both time and timescale ~~. Note~~and that, at a given time, events of different timescales can occur

simultaneously. ~~What one may subjectively interpret as a single event in the input time series is~~

330 ~~generally quantified by this definition as multiple coincident events at a variety of timescales~~

331 ~~each with a different power (e.g. Figure 1, panel c).~~

332 ~~We refer to contiguous regions of statistical significance (in time and timescale) as "event~~

333 ~~clusters" (note that no statistical clustering is performed).~~

334 3.1.3. Step 1c. Sample observed events to an event-set relevant to analysis

335 Step 1b results in the identification of all events at all timescales and times. In this sub-

336 step, the event space is sampled to suit the particular evaluation. Torrence and Compo (1998)

337 offer two methods to smoothing the wavelet plot that can increase significance and confidence:

338 (i) averaging in time (over timescale) or (ii) averaging in timescale (over time). Because the goal

339 of this paper is to evaluate model timing errors over long simulation periods, we choose to

340 sample the event space based on ~~dominant timescales in the time-averaged observed wavelet~~

341 ~~spectra~~averaging in timescale. Although for some locations there may be physical reasons to

342 expect certain timescales to be important (e.g., seasonal cycle of snowmelt), the most important

343 timescales at which hydrologic signals occur at a particular location are not necessarily known *a*

344 *priori*. Averaging events in timescale can provide a useful diagnostic ~~be used to~~by identify~~ing~~~~y~~

345 the ~~most~~ dominant, or ~~what we call~~ "characteristic", timescales for a given time series~~, which is~~

346 ~~useful for model diagnostics~~. Averaging many events in timescale can fil~~i~~ter noise and help

347 reveal the expected timescales of dominant variability corresponding to different processes or

348 sets of processes. ~~If one suspected nontstationarity in the timescales dominant variability over the~~

349 ~~timeseries, a different approach such as a moving average in timescale could be employed.~~~~The~~

350 ~~assumption is that identifiable sets of processes of interest are distinct in timescale, and that~~

351 ~~averaging over many events will reveal its expected value.~~

352　　　In our analysis we seek to uncover the dominant event timescales and to evaluate modeled

353　timing errors on these. The following bullets articulate our methodological choices for

354　summarizing observed events~~what was followed here~~~~For our application we choose to further~~

355　~~sub-sample the observed wavelet spectra by selecting, for each characteristic timescale, the most~~

356　~~powerful event within each event cluster. This is articulated in the following bullet~~s:

357　　　• *Calculate the average event power ~~across~~ in each timescale:* Considering only the

358　　　　statistically significant areas of the observed wavelet spectrum, calculate the average

359　　　　power ~~across~~ in each timescale (Figure 2c, right panel~~over time~~). We point out that

360　　　　calculating the average power over events is different than what is found by averaging

361　　　　across all time points, which doesn't take statistical significance into consideration

362　　　　(Figure 2b, right panel). ~~.~~

363

364　　　• *Identify timescales of absolute and local maxima in time--average power ~~maxima~~:* ~~By~~

365　　　　~~plotting the~~ After obtaining the average event ~~average~~ power ~~versus the~~ as a function

366　　　　timescale (Figure 2c, right panel), the local and absolute maximums for average event

367　　　　power can be determined. ~~(grey dots in Figure 2c, right); i~~In the Onion Creek case, there

368　　　　is a single maximum at 22 hours (grey dot in Figure 2c, right panel). ~~.~~ The timescales

369　　　　corresponding to the absolute and local maxima of the average power of the observed

370　　　　time series are called the characteristic timescales ~~of the observed wavelet spectrum~~used

371　　　　for evaluation. This is the first subset of events: all events that fall within the

372　　　　characteristic time-scales. For a single characteristic timescale, contiguous events in time

373　　　　are called event clusters (horizontal line in Figure 2d).

374

- *Identify events with maximum power ~~for~~ in each event cluster:* For all timescales, ~~As~~ ~~previously mentioned, events can also be grouped into "event clusters", that is, contiguous significant areas~~For each event cluster, ~~. We can use this to further sample from the event-set created in the last bullet: across each characteristic timescale,~~ we identify the event with maximum power in ~~for~~ each event cluster. This is the second event subset: all events with maximum power ~~for~~ in each cluster that fall~~s~~ within a characteristic timescale~~.~~ (star in Figure 2d); these are called cluster ~~maximums~~maxima (maxs).

*3.2. Step 2. Calculate Timing Errors*

Step 1 identifies ~~characteristics (cluster maxima) of~~ observed events by applying a wavelet transform to the observed time series. To calculate the timing error of a modeled time series, we perform its cross wavelet transform with the observed time series~~, as detailed in this section~~. Figure 3a shows the observed and modeled ~~time~~ time series used in our illustration of the methodology: ~~To illustrate Step 2, we use~~ the observed is the same isolated peak from Onion Creek, TX ~~(Figure 2a)~~, as in Figure 2a, and the synthetic modeled time series adds ~~and add~~ a prescribed timing error of +5 hours to the observed~~to every point in the original time series (Figure 3a) to create a synthetic time series~~. (Note that while the observed time series is identical in both, figures 2a and 3a have linear and log10 axes, respectively).

3.2.1. Step 2a. Apply cross-wavelet transform (XWT) to observations and simulations

~~For Step 2, we use the same Onion Creek, TX, peak from Figure 1a, and add a prescribed timing error of +5 hours to every point in the original time series (Figure 2a) to create a synthetic time series.~~The cross-wavelet transform (XWT) is performed between the observed and synthetic time series. ~~We perform the cross-wavelet transform between the observed and~~

17

399 ~~synthetic time series (Figure 2b).~~ Given the WT<u>s</u> of an observed time series $W_n^X(s)$ and a

400 modeled time series $W_n^Y(s)$, the cross-wavelet spectrum can be defined as:

401 $$W_n^{XY}(s) = W_n^X(s)W_n^{Y*}(s),$$

402 where the asterix ~~implies~~ <u>denotes</u> the complex conjugate. The cross-wavelet power is defined as

403 $|W_n^{XY}(s)|$ ~~.~~ <u>and signifies the joint power of the two time series. The XWT between the Onion</u>

404 <u>Creek observations and the synthetic 5 hour offset time series is</u>~~are~~ <u>shown in Figure 3b, with</u>

405 <u>power represented by the color scale.</u>

406 <u>Similar to Step 1b of the WT, we can also calculate areas of significance for the XWT</u>

407 <u>power as shown by the black contour in Figure 3b. For the XWT, significance is calculated with</u>

408 <u>respect to the theoretical background wavelet spectra of each time series (Torrence and Compo,</u>

409 <u>1998). We define XWT events as points of significant XWT power outside the COI. XWT</u>

410 <u>events indicate significant joint variability between the observed and modeled time series.</u>

411 <u>Below, in step 2d, we employ XWT events as a basis for identifying hits and misses on observed</u>

412 <u>events for which the timing errors are calculated. Figure 3c shows the intersection of the</u>

413 <u>observed events (colors) and the XWT events (dashed contour). As described later, this</u>

414 <u>intersection (inside dashed contour) is a region of hits where timing errors are considered valid.</u>

415 <u>Note that the early part of the observed events at shorter timescales is not in the XWT events.</u>

416 <u>This is because the timing offset in the modeled time series misses the early part of the observed</u>

417 <u>event in a way that depends on timescale.</u>

418 ~~Similar to Step 1b of the WT, we can also calculate the areas of significance for the~~

419 ~~XWTblack . In the next section, widentifyinghits and misses for the events for which the timing~~

420 ~~errors are calculatedthe confidence in theare quantified by looking at the percent hits, i.e., These~~

421 ~~are not the same as the areas of significance for the WT. The significant areas of the XWT vary~~

422 ~~with each simulation, and are therefore not useful for evaluation on their own. Nevertheless, we~~

423 ~~are interested in the overlap between the significant areas of the observed WT and the significant~~

424 ~~areas of the cross-wavelet transform, and this is used to quantify our confidence in the timing~~

425 ~~error estimate. We discuss this further in Step 2d.~~dashed

426 3.2.2. Step 2b. Calculate the cross-wavelet timing errors

427 For complex wavelets, such as the Morlet used in this paper, the individual WTs include

428 an imaginary component of the convolution. Together, the real and imaginary parts of the

429 convolution describe the phase of each time series with respect to the wavelet. The cross wavelet

430 transform combines the WTs in conjugate, allowing the calculation of a phase difference or

431 angle (radians) which can be computed as

432 ~~To calculate the timing errors, we first compute the phase angle of the cross-wavelet~~

433 ~~spectrum. The phase angle gives the phase difference and can be computed as~~:

434

435 $$\phi_n^{XY}(s) = tan^{-1}\left[\frac{\Im(\langle s^{-1}W_n^{XY}(s)\rangle)}{\Re(\langle s^{-1}W_n^{XY}(s)\rangle)}\right];$$

436

437 Where $\Im$ is the imaginary and $\Re$ is the real component of $W_n^{XY}(s)$. The arrows in Figure 3b

438 indicate the phase difference for our example case, which are used to calculate the timing errors.

439 Note that these are calculated at all points in the wavelet domain.

440 The distance around the phase circle at each timescale is the Fourier period (hours). We

441 convert the phase angle into the timing errors (hours) ~~We convert the phase angle into the timing~~

442 ~~error~~ as in Liu et al. (2011):

443

444 $$\Delta t_n^{XY}(s) = \phi_n^{XY}(s) * T/2\pi\frac{\cdot}{s}$$

445 where T is the equivalent Fourier period of the wavelet.

19

. Note that the maximum timing error which can be represented at each timescale is half the

Fourier period because the phase angle is in the interval (-pi, pi). In other words, only timescales

greater than 2E can accurately represent a timing error, E. Because the range of the *arctan*

function is limited by ±pi, true phase angles outside this range alias to angles inside this range.

(For example, the phase angles 1.05 * pi and -.95 * pi are both assigned to -.95*pi). Also note

that when the wavelet transforms are approximately antiphase, the computed phase differences

and timing errors produce corresponding bimodal distributions given noise in the data. Figure 3c

shows phase aliasing in the negative timing errors at timescales less than 10 hours, double the 5

hour synthetic timing error we introduced. The bimodality of the phase and timing are also seen

at ~~that~~ the 10hr timescale when the timing errors abruptly change sign (or phase by 2pi). We note

the convention used is that the XWT produces timing errors that are interpreted as "modeled

minus observed", i.e., ~~so that~~ positive values mean the model occurs after the observed. Positive

5 hour timing errors in Figure 3c describe that the model is "late" compared to the observations

as seen in the hydrographs in the top panel (a).

~~We convert the phase angle into the timing error as in Liu et al. (2011):~~

$$\Delta t_n^{XY}(s) = \phi_n^{XY}(s) * T/2\pi,$$

~~where T is the equivalent Fourier period of the wavelet.~~

~~The arrows in Figure 3b indicate the phase offset, which are used to calculate the timing~~

~~errors.~~

3.2.3. Step 2c. Subset cross-wavelet timing errors to sampled observed events

Step 2b results in an estimate of timing errors for all times and timescales in the cross-

wavelet transform space. In our application, we are interested in the timing errors that correspond

to the identified sample of *observed* events, especially for ~~events at the characteristic timescales~~

~~(the first event-set in step 1c) and for~~ the maximum power events in each cluster ~~(the second~~

470 ~~event set in step 1c). At~~ for each characteristic timescale~~, this provides a timing errors at each~~

471 ~~event cluster's maximum value~~. In the synthetic Onion Creek example, the point~~s~~ of interest in

472 the wavelet transform of the observed timeseries, used to sample the timing errors produced by

473 the XWT, ~~are~~is shown by the grey star in Figure 3c. ~~The latter provides a single timing error for~~

474 ~~each event cluster max at each characteristic timescale, which could be used in a post-processing~~

475 ~~step to provide a cluster-by-cluster timing correction, if desired.~~

476 ~~In Figure 3c, the timing error estimates show that for timescales greater than 10 hours, we~~

477 ~~get back the prescribed timing error of 5 hours, i.e., the scale must be at least double the timing~~

478 ~~error.~~ The results for the synthetic Onion Creek example are summarized in Table 2.~~;~~ F~~f~~or the

479 identified characteristic timescale of 22 hours in the observed wavelet power~~,~~ (which had an

480 average WT power of 555,700~~676598,000~~ m^6/s^2 ~~-~~ (~~from~~Figure 2c right), there was 1 event

481 cluster, and the timing error ~~for~~at the cluster max~~imum~~ was 5 hours (and it~~, which~~ occurred at

482 hour 37 of the time series). ~~maximum~~

483 3.2.4. Step 2d. ~~Quantify~~Filter Misses~~ Percent Hits~~

484 The premise of computing a timing error between the observed and modeled time series

485 is that they share common events which can be meaningfully compared. In a two-way

486 contingency analysis of events, a "hit" refers to when the modeled time series reproduces an

487 observed event. When the modeled time series fails to reproduce an observed event, it is termed

488 a "miss". In the case of a miss, it does not make sense to include the timing error in the overall

489 assessment. ~~Because the timing errors are calculated from the XWT, we choose to diagnose hits~~

490 ~~and misses based on the significance of the XWT.~~ Once the characteristic timescales of the

491 observed event spectrum are identified and event cluster maxima are located, timing errors are

492 obtained at these locations in the XWT. In this step, the significance of the XWT on these event

493 cluster maxima is used to decide if the model produced a hit or a miss for each point and to

494 determine if the timing error is valid. ~~For a single cluster max, such as shown in Figure 3c, the~~

495 ~~XWT significance is either True or False, the point is either a hit or a miss.~~ As previewed above,

496 Figure 3c shows ~~shows~~ the observed events (colors) and the XWT events (dashed contour).

497 Regions of intersection between observed events and XWT events are considered model hits~~the~~

498 ~~intersection of the observed events (colors) and the XWT events (dashed contour)~~ and ~~.~~observed

499 events falling outside the XWT events are considered misses. Because we constrain our analysis

500 to observed events in the wavelet power spectrum, we do not consider either of the remaining

501 categories in a 2-way analysis (false alarms and correct negatives). We note that a complete 2-

502 way event analysis could alternatively be constructed in the wavelet domain based on the Venn

503 diagram of the observed and modeled events without necessarily using the XWT. W~~As~~

504 ~~mentioned, w~~e choose to use the XWT events because the XWT is the basis of the timing errors.

505     In the synthetic example of Onion Creek, ~~This is a region of hits. For~~ a single

506 characteristic timescale and event cluster yields a single cluster max~~, such as~~ as shown by the

507 star in Figure 3c. Because this star falls both within the observed and XWT events, it is a hit and

508 the timing error at that point is valid (Table 2)~~, the XWT significance is either True or False, the~~

509 ~~cluster max (star)~~ point is either a hit or a miss. ~~Table 2 summarizes the results of the timing~~

510 ~~error analysis for this synthetic example. We can see the prescribed 5 hour offset is recovered by~~

511 ~~the calculation and that the timing error is valid because the observed event was reproduced by~~

512 ~~the model (a hit).~~ For a longer time series, as seen in subsequent examples, a useful diagnostic

513 and compliment to timing error statistics at each characteristic timescale is the percent hits.

514 When summarizing timing errors statistics for a timescale, we drop misses from the calculation

515 and the % hits indicates what portion of the time series was dropped (% misses = 100 - % hits).

516 In our tables we provided timing error statistics only for hits~~this way as well as over all observed~~

517 ~~events to reveal the impact of dropping misses~~.

518

519 ~~It is important to point out that for other applications, there could be other ways to~~

520 ~~interrogate the timing errors that result from the cross-wavelet transform. Some of these~~

521 ~~possibilities are noted in the Discussion section.~~

522 ~~3.2.4. Step 2d. Quantify Percent Hits~~ Quantify the confidence in the timing error estimate

523 ~~we report the hits to provide an assessment of the confidence in the timing error~~

524 ~~assessment~~To interpret our confidence in the timing error estimate, we can examine ~~if the cluster~~

525 ~~maxs are~~ the overlap between the significant areas of the observed WT and the significant areas

526 of the XWT.

527 ~~We can look at percent (%) overlap, that is, how many of the XWT events overlap with the WT~~

528 ~~events, either for all events or for the sampled event-sets. An overlap close to 0% would indicate~~

529 ~~that the model did not do a good job of simulating the observations – or it is a "miss" (flood is~~

530 ~~observed but not forecasted). If the overlap was 100%, it would be close to a perfect simulation.~~

531 ~~Second, I~~if we are looking at a single timing error for each event cluster, we may look to see if

532 ~~that event is significant in the XWT. If it is not, it gives us less confidence in the estimate.~~ In

533 ~~Table 2, we can see for the prescribed 5 hour offset example, the cluster max was significant in~~

534 ~~the XWT. When there are multiple clusters for a given characteristic timescale, the %~~

535 ~~significance can be calculated as the ratio of the # of significant cluster max's to the total number~~

536 ~~of cluster maxs.~~

537 ~~We note that because we are calculating timing errors in terms of observed events, there is no~~

538 ~~information about "false alarms", where a flood is forecasted but not observed.~~

## 4. Application of the Framework

The methodology developed in this paper is implemented in the R language and is made publicly available, as detailed in the code availability section at the end of the manuscript.

### 4.1. Data

The application of the methodology is illustrated using real and simulated stream discharge (streamflow, m3/s) data from four U.S. Geological Survey (USGS) stream gage locations: Onion Creek at US Highway 183, Austin, Texas (Onion Creek, TX; USGS site number 08159000), Taylor River at Taylor Park, Colorado (Taylor River, CO; USGS site number 09107000), Pemigewasset River at Woodstock, New Hampshire (Pemigewasset River, NH; USGS site number 01075000), and Bad River near Fort Pierre, South Dakota (Bad River, SD; USGS site number 06441500). We use the USGS instantaneous observations averaged on an hourly basis.

NOAA's National Water Model (NWM, https://www.nco.ncep.noaa.gov/pmb/products/nwm/) is an operational model that produces hydrologic analyses and forecasts over the continental United States (CONUS) and Hawaii (as of version 2.0). The model is forced by downscaled atmospheric states and fluxes from NOAA's operational weather models. Next, the NoahMP (Niu et al 2011) land surface model calculates energy and water states and fluxes. Water fluxes propagate down the model chain through overland and subsurface (soil and aquifer representations) water routing schemes to reach a stream channel model. The NWM applies the three parameter Muskingum-Cunge river routing scheme to a modified version of the NHD-Plus version 2 (McKay et al. 2012) river network representation.

In this study, NWM simulations are taken from each version's retrospective runs (https://docs.opendata.aws/nwm-archive/readme.html). These are continuous simulations (not cycles) run for the period October 2010 to November 2016 and forced by the National Data Assimilation System (NLDAS)-2 product as atmospheric conditions. The nudging data

563 assimilation was not applied in these runs either. We use NWM discharge simulations from

564 versions V1.0, V1.1, and V1.2 (not all version may be publicly available).

565 To apply the methodology, we note that the observed and simulated datasets must be paired

566 (overlapping). Further, for evaluation, any new simulation must also be paired with the observed.

567 Missing data, which is common in observed time series, can be problematic and can result in

568 false significance. We account for this our methodology by calculating the XT and XWT on each

569 complete time series. This will be illustrated in the forthcoming example at Taylor River, CO.

570 *4.2. Application*

571 For illustration purposes we apply Steps 1 and 2 to an observed time series in Onion Creek, TX;

572 for simplicity, we select an isolated peak (Figure 1a). First, we apply the wavelet transform to the

573 observations (Figure 1b). This shows the time series in terms of its power by time and timescale,

574 with warmer colors indicating more power. The black outline shows the areas of significance and

575 the muted colors indicate the COI. To determine all observed events, we identify all the points

576 that are significant and outside the COI (Figure 1c). Next, we average the power across each

577 timescale: to the right of Figure 1b we show power averaged across all points for each timescale,

578 and to the right of Figure 1c we show power averaged across just the events for each timescale.

579 The latter is the one used to identify our characteristic scales. In this case, there is a single

580 maximum at 22 hours. For the characteristic timescale, we see there is only 1 event cluster and

581 the event with maximum power is marked with a star (Figure 1d).

582 For Step 2, we use the same Onion Creek, TX, peak from Figure 1a, and add a prescribed

583 timing error of +5 hours to every point in the original time series (Figure 2a) to create a synthetic

584 time series. We perform the cross-wavelet transform between the observed and synthetic time

585 series (Figure 2b). The arrows in Figure 2b indicate the phase offset, which are used to calculate

586 the timing error (Figure 2c). The timing error estimates show that for timescales greater than 10

587　~~hours, we get back the prescribed timing error of 5 hours, i.e., the scale must be at least double~~

588　~~the timing error. In this case, because we are adding a prescribed error, the error is approximately~~

589　~~5 hours for all events, including for the characteristic timescale of 22 hours.~~

590　~~Finally, we repeat Step 2, but compare the observation of this event to actual model data~~

591　~~from NWM V1.2. This shows that the model is early (Supplemental Figure 2a). We perform the~~

592　~~cross wavelet transform (Supplemental Figure 2b) and examine the timing error (Supplemental~~

593　~~Figure 2c). Table 1 summarizes the results: the mean error across the 22-hour characteristic~~

594　~~timescale is -3.2 hours, as is the error for the cluster's maximum power. All events in the cluster~~

595　~~are also significant in the XWT (100%), and the cluster maximum is also significant, providing~~

596　~~confidence in this timing error estimation.~~

597

598　**45. Results**

599　In the previous section, we illustrate the method using an isolated peak and a prescribed

600　timing error. In this section, we ~~further~~ demonstrate the method~~, increasing the complexity by~~

601　using ~~~~NWM model ~~modeled~~ simulat~~ed data~~ions which introduce greater complexity ~~~~and longer

602　time series ~~is used from several locations and time series to highlight the features of the method.~~

603　Finally, we show ~~, finishing with~~ version-over-version comparisons for 5-year simulations to

604　illustrate the utility for evaluation.

605　4.1 Demonstration using NWM data

606　*~~5.1.~~Pemigewasset River, NH*

607　This example uses a three-month time series from the Pemigewasset River, NH, to ~~. First,~~

608　~~we~~ examine ~~a three-month time series that exhibits~~ multiple peaks ~~above a base flow~~in the

609　hydrograph (Figure 34a). By eye, it is fairly straightforward to pick out three main peaks. From

610   Step 1 of our method, applying t~~T~~he wavelet transform on~~f~~ the observations (Figure ~~3~~4b and

611   ~~3~~4c)~~,~~. reveals up to three event clusters, depending on the characteristic timescale examined

612   (Figure 4~~3~~d). When we plot the average event power by timescale (right of Figure ~~3~~4c), we see

613   that there are nine relative maxima (small grey dots) – hence there are 9 characteristic scales for

614   this example. ~~~~The ~~~~cluster maxima (grey stars) for each observed event cluster are shown in

615   Figure 4d.

616      In Step 2, we compare the ~~same~~ observed time series ~~from step 1~~ with ~~output~~ the simulation

617   from the NWM V1.2 (Figure 4~~5~~a): a)~~,~~ apply the cross-wavelet transform (Figure 4~~5~~b colors), b)

618   ~~and~~ calculate the timing error for all observed events (Figure 5b arrows), c) ~~~~subset the timing

619   errors to the observed cluster maxima (Figure 5c stars), and d) retain only modeled hits (Figure

620   5c stars within the dashed contours). ~~(Figure 4~~5~~c).  As previously mentioned, we are interested in~~

621   ~~the timing errors corresponding to observed events at the characteristic timescales. In~~Table 3

622   ~~Figure 5a, the panels are~~is ordered by characteristic timescales from highest to lowest average

623   power~~-~~; we only show the top 5 characteristic scales~~, using the first subset of events, grouped by~~

624   ~~cluster. The first panel, where timescale = 24.8 hours, is t~~The absolute maximum of the time

625   average event spectrum has a timescale = 24.8 hours; ~~(Figure 4c). This shows two cluster~~

626   ~~distributions:~~ for cluster one, the model is late~~close to on-time (-0.052 hr)~~ nearly 11 hours late

627   and cluster two is late~~early (7hr~~3.5 hours), both are hits, ~~so~~and the average timing error is 3.5

628   hours late. However, for the next timescale (=27.8 hr), ~~one of the~~the third cluster maximum~~s~~ is a

629   miss, so its ~~the~~ timing error is reported as a NA, and is not included in the average. This miss can

630   be seen in Figure 5c where the last star falls just outside the XWT events. Moreover, this miss

631   can also be interpreted from the comparison of the hydrographs in Figure 5a where the modeled

632   third peak does not reasonably approximate the magnitude of the observed peak. Interestingly,

633 the miss is a narrow miss at the shorter timescale of 27.8 hours while the associated (3$^{rd}$) cluster

634 maxima at the next most powerful characteristic timescale (33.1 hours) is a hit. This reflects that

635 hydrograph is insufficiently peaked for this event but does have some of the observed, lower-

636 frequency variability. Overall, this next most important characteristic timescale of 33.1 hours has

637 timing results similar to the 27.8 hour timescale with the exception of the third cluster maximum.

638 This raises the question if these are distinct characteristic timescales. In the Discussion and

639 Conclusions section we discuss ~~We point out that the~~ for most events, and cluster two shows the

640 ~~model is early; the dark shading indicates that most of the events are significant in the XWT. The~~

641 ~~next two dominant scales of~~ 27.8 and 33.1 hours ~~have similar average power and are of the same~~

642 ~~order of magnitude at 27.8 hours and 33.1 hours; if we had applied~~ smoothing ~~to the graph of~~ the

643 time average event ~~average~~ power by timescale to address this issue~~, these relative maxima~~

644 ~~would smooth out. We will revisit this in the Discussion, when we discuss pathways to~~

645 ~~implementation~~ in the Discussion and Conclusions.

646     The characteristic timescale with the ~~next~~ 4$^{th}$ highest time-average power ~~maxima~~ occurs at

647 111 hours, which is a different order of magnitude, suggesting that this may have a different

648 physical process driving it. ~~This~~ At this timescale, the ~~shows the~~ model is ~~to be~~ late in~~for~~ both

649 event clusters (10 and 16 hours). ~~, and r~~Results are similar for ~~a~~the next timescale of 148 hours.

650 We don't show results for the remaining 4 characteristic time scales with lower average power,

651 since they have similar characteristic timescale values and associated timing errors to what has

652 already been shown.

653 ~~We can see how looking at the timing errors using the cluster distributions will get harder as the~~

654 ~~number of clusters increase, so it is also useful to summarize the information by looking at each~~

655 ~~cluster mean and max. If we run the methodology on the full 5-year Pemigewasset River time~~

28

series, we can compare the mean and max timing errors for each characteristic time scale using box plots where the outline is shaded by the average confidence (Supplemental Figure 3). Table 2 summarizes this information. For example, the absolute maxima, at the 17.5 hour timescale has 86 clusters, and a timing error centered around zero (-0.43 hours), 75% of which are significant in the XWT. This is very similar to the results for the cluster max, as it is for the rest of the characteristic time scales. One other thing to note is that as expected, because the characteristic time scales are data driven, they are not the same as they were for the 3-month period.

5.2. Bad River, SD

The second example uses a two-month time series from the Bad River, SD, to illustrate the concept of consecutive peaks (Figure 6a). Whereas in the previous example it was fairly straightforward to pick out 3 distinct peaks, in this time series, there is one noticeable peak centered around June the 1st, with smaller peaks preceding and following it. The question is whether or not this is one event cluster or multiple? Looking at the wavelet transform (Figure 6b and 6c), we can see that for smaller timescales, there are more clusters, but for longer timescales, they are considered a single cluster.

In Step 2, we compare the same time series with output from NWM V1.2 (Figure 7a), calculate the cross-wavelet transform (Figure 7b), and calculate the timing error (Figure 7c). The timing error figure shows a sign switch: for longer timescales (i.e., when the peaks are considered part of a single event cluster), the model is early, but for shorter time scales (i.e., when the peaks are each considered their own cluster), the model is late. This is an important point: corrections at one scale may worsen timing error (or other metrics) at other scales.

This example has another interesting feature: namely that there is a false alarm in the model just before July 15. We note that because of our methodology, there is no observed event at that

679   ~~time, and therefore no timing error to be calculated, that is there is no information in the timing~~

680   ~~error statistics in terms of false alarms.~~

681   ~~5.3.~~ *Taylor River, CO*

682

683       In this example, we ~~will~~ examine a <u>one-year</u> time series from Taylor River, CO, that

684   illustrates <u>hydrograph</u> peaks ~~that are~~ driven by different processes. The Taylor River is in a

685   mountainous area where the spring hydrology is dominated by snowmelt runoff. ~~To start, we will~~

686   ~~look at a portion of the spring melt season, where we can visibly see a diurnal signal (Figure 8).~~

687   ~~However, while it's easy to see that the model is too high in amplitude, it's hard to visually tell~~

688   ~~much about the timing error. Figure 9 shows that for the characteristic time scale of 23.4 hours,~~

689   ~~the model is usually early, with high confidence.~~

690       ~~Supplemental Figure 4a~~<u>Figure 6</u>~~7~~<u>a</u> shows ~~a~~<u>the</u> ~~year-long~~ time series from Taylor River, CO,

691   where we can see the snowmelt runoff in spring <u>,</u> ~~but~~<u>and</u> also several peaks in summer, likely

692   driven by summer rains. ~~Supplemental Figure 4~~<u>Figure 6</u>~~7~~<u>b</u> shows the WT, and ~~also~~ illustrates

693   <u>how missing data is handled: this results in additional COIs (muted colors) to account for the</u>

694   <u>edge effects, and areas of the COI are ignored in our analyses.</u>

695       From the <u>statistically significant events in the</u> WT, we ~~again~~ see the peak in the characteristic

696   time scales at about 24~~4~~ hours (right of ~~Supplemental~~ Figure 6~~7~~<u>4</u>c), ~~but~~ <u>and</u> there is another

697   maxima at 99 and 118 hour timescales, relating to flows from the summer rains. ~~This non-~~

698   ~~stationarity dominant timescale is evident in the wavelet power (Figure 6b and 6c)~~. <u>In Step 2, we</u>

699   <u>compare the</u> ~~same~~<u>observed</u> <u>time series with</u> ~~output~~<u>the simulation from the NWM V1.2 (Figure</u>

700   <u>7</u>~~8~~<u>a); here it is useful to zoom into the spring melt season time series (Figure 8</u>~~9~~<u>), where we see</u>

701   <u>that the amplitude of the diurnal signal is too high, but it's hard to visually tell much about the</u>

702   <u>timing error. Next, the cross-wavelet transform (Figure 7</u>~~8~~<u>b) and timing errors are calculated</u>

703 (Figure ~~8~~7c). The results are summarized in Table 4. ~~Looking at Figure 10, s~~Starting with the

704 dominant 2~~4~~3 hour timescale, we see that ~~for the~~there are 11 clusters, that 73% (=8/11 cluster

705 maxima) are hits, and that the model is ~~which are~~ generally early (the mean is ~~4.~~6 hours early)~~,~~

706 ~~and that 73% (=8/11 cluster maxs) are significant in the XWT.~~ ~~that are significant in the XWT,~~

707 ~~the model is generally early.~~ For the 118 and 99 hour timescale~~, the model is also early, b~~s, ~~u~~t

708 ~~none of the those cluster events~~maxs are ~~are not statistically significant in the XWT~~there are no

709 hits. ~~(0%).~~ This suggests that we are confident in the ~~early~~ timing errors of the model for the

710 diurnal snowmelt cycle, and these timing errors can ~~this could~~ be used as ~~qualitative~~ guidance for

711 model performance and model improvements~~at this site until the model performance is~~

712 ~~improved~~. However, the model does not successfully reproduce key variability during the

713 summer and timing errors are not valid at this timescale~~can not be used to evaluate or guide~~

714 ~~model improvements during this time~~. ~~we show that it is less reliable for the early timing errors~~

715 ~~for the summer peaks.~~ This underscores the key point that timing errors are timescale dependent,

716 and can help diagnose which processes to target for improvements.

717 ~~Supplemental Figure 4b also illustrates how missing data is handled: this results in additional~~

718 ~~COIs (muted colors) to account for the edge effects, and areas of the COI are ignored in our~~

719 ~~analyses.~~

720

721 4.2 Evaluating Model Performance
722
723 Finally, we show how the methodology can be used for evaluating performance changes

724 across NWM versions. We point out that none of the NWM version upgrades were targeting

725 timing errors, so these results just provide a demonstration. We use ~~a~~ 5-year observed and

726 modeled ~~overlapping time series and cluster max for the results~~time series at the three locations~~,~~:

727 Onion Creek, TX, and Pemigewasset River, NH, but cluster mean results were similar (not

728 shown), and Taylor River, CO..

729 For Onion Creek, Table 5 summarizes the results for the three most important timescales and

730 Figure 9 provides a graphical representation of these timing errors (hits only). For the NWM

731 V1.0 for Onion Creek, we see that fFor the dominant 29.5 hour timescale and for all model

732 versions, there were 197 cluster maximas, all of which were hits89.5% of which were hits, with a

733 median timing error of 1.4 hours early, for which the median timing error is -1.4 hours, and all

734 were significant in the XWT (Table 35). However, the model showeds progressively earlier

735 timing errors with increasing version (Figure 9). The results are similar for the other two

736 characteristic timescales.

737 Comparing V1.0, V1.1, and V1.2, the results for Onion Creek show that the median timing

738 error has gotten slightly earlier (worse), although the distribution became tighter from V1.0 to

739 V1.1 and V1.2 (Figure 1011). In Figure 1011, the dark blue color of the boxplot outline indicates

740 that there is high confidence in the timing error, as the overlapping significance is close to 100%

741 for the top three characteristic timescales. Using the 5-year overlapping time series forFor

742 Pemigewasset River, NH, Table 6 summarizes the results for the 3 most important timescales

743 and Figure 10 provides a graphical representation of the timing errors (hits only). At this

744 location, we see that the median timing error improved by improved with NWM V1.2, getting

745 closer to zero,. While the distribution of the timing errors became less biased than the previous

746 versions, it also became but that the distribution became wider (Figure 1012). Over the

747 timeseries, there were between 59 and 76 event clusters. Interestingly, the hit rate for all

748 timescales was best for NWM V1.1 though its timing errors are broadly the worst. Again, the

749 confidence is fairly highhits are fairly high (>80%) across characteristic time scales and versions

~~(Table 46)~~From NWM V1.0 to NWM V1.2, improvements to both hit rate and median timing errors were obtained at all timescales.

~~, and >60 clusters were used in the estimations. Using 5-years from~~For Taylor River, ~~CO~~ ~~(Supplemental~~ Table 7~~2, Supplemental Figure 5),~~ summarizes the results for the 2 most important timescales. ~~we see that f~~For the characteristic timescale of 235 hours (~10 days) there are only 4 event clusters, ~~has low confidencebut there are not~~ and each model version has only 1 hit. The timing of this hit improves by roughly half its error from NWM V1.0 to NWM V1.2 in going from 16 to 9 hours. ~~many hits (~25%) for the 4 sampled clusters; T~~the ~~timescale of~~ 23.4 hour timescale has 41 event clusters with a hit rate varying considerably by version. ~~s has a~~ The median timing error ~~that is~~ is fairly consistent with version, however, ranging from 6 to 7 hours ~~ly~~ early ~~by around 6 hours, with the version model confidencehits~~ ranging from 44% to 67% (Supplemental Table 2~~7~~). ~~Results for the Bad River can be seen in Supplemental Table 3 and Supplemental Figure 6.~~

**6. Discussion and Conclusions**

In this paper, we develop a systematic, data-driven methodology to objectively identify timeseries (hydrograph) events and estimate timing errors in large-sample, high-resolution hydrologic models. The method was developed towards several intended uses: Primarily, it was developed for model evaluation, so that model performance can be documented in terms of defined standards. We illustrate this with the version-over-version NWM comparisons. Second, it can be used for model development, whereby potential timing error sources can be diagnosed (by timescale) and targeted for improvement. Related to this point, given the advantages of

774 calibrating using multiple-criteria (e.g., Gupta et al. 1998), timing errors could be used as part of

775 a larger calibration strategy. However, ~~as noted in the consecutive peaks example for the Bad~~

776 ~~River,~~ minimizing timing errors at one timescale may not translate to improvements in timing

777 errors (or other metrics) at other timescales. ~~-~~Wavelet analysis has also been used directly as an

778 objective function for calibration, although a difficulty is in determining the similarity measure

779 to use (e.g. Schaefli and Zehe 2009, Rathinasamy et al. 2014). Future research will investigate

780 the ~~properties~~ application of~~of the~~ timing errors presented here for calibration purposes. Finally,

781 the approach can be used for model interpretation and forecast guidance~~,~~ as estimating timing

782 errors provides ~~a~~ characterization of the timing uncertainty (i.e., for a given timescale, the model

783 is generally late or early)~~, as well as a measure of the~~ or confidence~~, that could be useful for~~

784 ~~qualitative forecast guidance~~.

785     Given the fact that several subjective choices were made specific to our application and

786 goals, ~~we think~~ it is important to highlight that we have made the analysis framework openly

787 available (detailed in the code availability section below), so the method can be adapted,

788 extended, or refined by the community right away. ~~For instance, because of our focus on model~~

789 ~~evaluation and development, we use the observed WT to identify events. However, i~~In other

790 ~~instances~~ applications it might be sufficient to only sample events that are in the significant areas

791 ~~of the XWT (essentially to identify the characteristic scales and event-set directly from the XWT~~

792 ~~instead of from the WT). This might be reasonable for applications that are more focused on~~

793 ~~model interpretation in a real-time forecasting mode, but it would not allow for version~~

794 ~~comparison and it is not guaranteed that all the important characteristic scales would be~~

795 ~~identified (i.e., the model may not capture some real-world processes, and therefore miss the~~

796 ~~associated characteristic timescales).~~ We ~~only~~ look at ~~the~~ timing errors from an observed event-

34

797 set relevant to our analysis, but there are other ways to subset the events that might be more

798 suitable to other applications. For ~~instance~~example, we focus on the event cluster maxima, but

799 one could also examine the event cluster means or the local maxima along time. ~~Also~~Another

800 alternative to~~, instead of~~ finding the ~~event of maximum power in each~~event cluster maxima (i.e.,

801 for a given timescale)~~, it~~ would be ~~possible ~~to identify the event with maximum power in

802 "islands of significance" across timescales, i.e., contiguous regions of ~~significant areas~~

803 contiguous significance ~~in time ~~across both time and timescale~~s. However, t~~This approach would

804 ignore ~~s ~~that multiple frequencies can be important at once. Also~~-~~ ~~and ~~defining ~~the ~~such islands is

805 ~~also ~~not straightforward~~when connected is problematic~~. ~~Yet another~~A differen~~ce~~t approach could

806 be desirable ~~I~~if one suspected non-stationarity in the characteristic timescales over the

807 ~~timeseries~~time series~~,~~. ~~t~~Then perhaps ~~a different approach such as ~~a moving average in timescale

808 could be employed to identify characteristic timescales. ~~For instance, Further, ~~In our approach,

809 we define the event set broadly~~,~~. However, ~~but ~~it could be subset ~~for high peak or ~~using

810 streamflow thresholds (e.g. for flooding events) to compare events in the wavelet domain with

811 traditional peak-over-threshold ~~approaches~~events. For example, ~~Supplemental ~~Figure 11~~7~~ shows

812 the maximum streamflows for the event-set from the 5 year run at Taylor River. This figure

813 shows that all events identified by the algorithm are not necessarily high flow events (i.e., the

814 maximum streamflow peaks are lower for the 23.4 hour timescale as compared to the 235.6 hour

815 timescale). To compare with traditional peak-over-threshold approaches, this event-set could be

816 filtered to include only events above a given threshold (i.e., events in both the wavelet and time

817 domains). ~~; this event-set could be filtered to include only events above a given ~~~~streamflow~~

818 ~~threshold (i.e. events in both the wavelet and time domains). ~~~~The method provides a~~

819 ~~quantification of the confidence in the~~~~percent hits for the timing errors; however, ~~~~, and we~~

820 ~~include all timing errors in our summaries, whether they are hits or misses. However, i~~It might

821 make more sense to drop ~~those points~~ in the timing error assessment,  that do not have a high

822 ~~confidence (i.e., with a low percent of events that significantly overlap between the XT and the~~

823 ~~XWT) and to~~ and to only calculate the timing errors on hits.~~flag those events as misses.~~

824 Another point that arises is how many characteristic timescales should be examined and

825 the similarity of adjacent characteristic timescales. ~~Here~~In our method, we average the power

826 ~~across~~ in timescales ~~-~~and identify characteristic scales to be at every absolute and relative

827 maxima. As seen in the illustrative examples, this can result in multiple characteristic scales,

828 some of which can be quite similar, suggesting that events at those scales are from similar or

829 related processes. One solution could be to smooth the average power by timescale, which would

830 reduce the number of local maxima, or to look at timing errors within a band of timescales. It is

831 also important to note that the characteristic scales are data-driven, so they will change with

832 different lengths of observed time series. Longer runs capture more events and should converge

833 on the more dominant timescales and events for a location. However, for performance

834 evaluation, overlapping time periods for observed and modeled time series are needed.

835 In our application of the WT, we follow Liu et al. (2011) and select the Morlet as the

836 mother wavelet. However, results are sensitive to the mother wavelet selected. Further discussion

837 of mother wavelet choices can be found in Torrence and Compo (1998) and in ElSaadani and

838 Krajewski (2017).

839 In ~~short~~summary, this paper provides a systematic, flexible, and computationally efficient

840 methodology for calculating model timing errors that is appropriate for model evaluation and

841 comparison, and is useful for model development and guidance. Based on the wavelet transform,

842 the method introduces timescale as a property of timing errors. The approach also identifies

843 streamflow events in the observed and modeled timeseries and only evaluates timing errors for

844 modeled events which are hits in a 2-way contingency analysis. Future work will apply the

845 approach to identify characteristic timescales across the United States, as well as to assess the

846 associated timing errors in the NWM.

847 **Code/Data Availability**

848 The code for reproducing the figures and tables in this paper are provided in the public github

849 repository https://github.com/NCAR/wavelet_timing with instructions for installing

850 dependencies. The core code used in the above repository is provided in the public "rwrfhydro"

851 R package https://github.com/NCAR/rwrfhydro. The code is written in the open-source R

852 language (R Core Team 2019) and builds off multiple, existing R packages. Most notably the

853 wavelet and cross-wavelet analyses are performed using the "biwavelet" package (Gouhier et al.

854 2018).

855 We emphasize that the analysis framework is meant to be flexible and adapted to similar

856 applications where different statistics may be desired. The figures created are specific to the

857 applications in this paper but provide a starting point for other work.

858 The code for reproducing the figures in this paper as well as extended vignettes/notebooks are

859 provided in public github repository https://github.com/NCAR/wavelet_timing. In addition to

860 reproducing the analyses and figures in this paper, several jupyter notebooks provide more

861 detailed analyses of the time series included in this paper. We emphasize that the analysis

862 framework is meant to be flexible and adapted to similar applications where different statistics

863 may be desired. The figures created are specific to the applications in this paper but provide a

864 starting point for other work.

37

865 ~~The core code is provided in the public "rwrfhydro" R package~~

866 ~~https://github.com/NCAR/rwrfhydro. The package can be installed as described by the~~

867 ~~README document in the repository and in the Supplemental Online Materials for this paper.~~

868 ~~The code is written in the open-source R language (R Core Team 2019) and builds off multiple,~~

869 ~~existing R packages. Most notably the wavelet and cross-wavelet analyses are performed using~~

870 ~~the "biwavelet" package (Gouhier et al. 2018).~~

871 **Credit Author Statement**

872 ET and JLM collaborated to develop the methodology. ET led the results analysis and

873 manuscript preparation and revisions. JLM developed the initial idea for the work, the open

874 source software, and visualizations.

875 **Competing interests**. The authors declare that they have no conflict of interest.
876

886
887 **References**:
888
889 Bogner, K., Pappenberger, F., 2011. Multiscale error analysis, correction, and predictive
890 uncertainty estimation in a flood forecasting system. Water Resour. Res. 47, W07524,
891 doi:10.1029/2010WR009137.
892

893     Coles, S., 2001. An Introduction to Statistical Modeling of Extreme Values, Springer Ser. Stat.
894     Springer, London.
895

896     Daubechies, I., 1990. The wavelet transform time-frequency localization and signal analysis.
897     IEEE Trans. Inform. Theory 36, 961-1004.
898

899     ElSaadani, M., Krajewski, W. F., 2017. A time-based framework for evaluating hydrologic
900     routing methodologies using wavelet transform. *Journal of Water Resource and Protection*, *9*(7),
901     723–744.
902

903     Ehret, U., Zehe, E., 2011. Series distance-an intuitive metric to quantify hydrograph similarity in
904     terms of occurrence, amplitude and timing of hydrological events. *Hydrology and Earth System*
905     *Sciences*, *15*, 877–896. https://doi.org/10.5194/hess-15-877-2011
906

907     Gochis, D.J., M. Barlage, R. Cabell, M. Casali, A. Dugger, K. FitzGerald, M. McAllister, J.
908     McCreight, A. RafieeiNasab, L. Read, K. Sampson, D. Yates, Y. Zhang (2020). The WRF-
909     Hydro® modeling system technical description, (Version 5.1.1). *NCAR Technical Note.* 107
910     pages. Available online at:
911     https://ral.ucar.edu/sites/default/files/public/WRFHydroV511TechnicalDescription.pdf. Source
912     Code DOI:10.5281/zenodo.3625238
913

914     Gouhier, T.C., Grinsted, A., Simko, V., 2018. R package biwavelet: Conduct Univariate and
915     Bivariate Wavelet Analyses (Version 0.20.17), https://github.com/tgouhier/biwavelet.
916     Gupta, H.V., Kling, H., Yilmaz, K.K., Martinez, G. F., 2009. Decomposition of the mean
917     squared error and NSE performance criteria: Implications for improving hydrological modelling.
918     Journal of hydrology, 377(1-2), 80-91.
919

920     Gupta, H.V., Perrin, C., Blöschl, G., Montanari, A., Kumar, R., Clark, M., Andréassian, V.,
921     2014. Large-sample hydrology: a need to balance depth with breadth. Hydrol. Earth Syst. Sci.
922     18, 463-477, https://doi.org/10.5194/hess-18-463-2014.
923

924     Gupta, H.V., Sorooshian, S., Yapo, P.O., 1998. Towards improved calibration of hydrologic
925     models: multiple and non-commensurable measures of information. Water Resources Research
926     34(4): 751–763.
927

928     Gupta, H.V., Wagener, T., Liu, Y., 2008. Reconciling theory with observations: elements of a
929     diagnostic approach to model evaluation. *Hydrological Processes*, *22*(March), 3802–3813.
930     https://doi.org/10.1002/hyp
931

932     Koskelo, A. I., Fisher, T. R., Utz, R. M., and Jordan, T. E., 2012. A new precipitation-based
933     method of baseflow separation and event identification for small watersheds (< 50 km2), J.
934     Hydrol., 450–451, 267–278, https://doi.org/10.1016/j.jhydrol.2012.04.055.
935

936     Lane, S.N., 2007. Assessment of rainfall–runoff models based upon wavelet analysis. Hydrol.
937     Processes 21, 586–607, https://doi.org/10.1002/hyp.6249.
938

939 Liu, Y., Liang, X.S., Weisberg, R.H., 2007. Rectification of the bias in the wavelet power
940 spectrum. J. Atmos. Oceanic Technol. 24, 2093–2102.
941
942 Liu, Y., Brown, J., Demargne, J., Seo, D. J., 2011. A wavelet-based approach to assessing timing
943 errors in hydrologic predictions. J. Hydrol. 397(3–4), 210–224.
944 http://doi.org/10.1016/j.jhydrol.2010.11.040
945
946 Luo, Y. Q., Randerson, J.T., Abramowitz, G., Bacour, C., Blyth, E., Carvalhais, N., … Orme, C.
947 E., 2012. A framework for benchmarking land models. Biogeosciences 9, 3857-3874,
948 https://doi.org/10.5194/bg-9-3857-2012.
949
950 McKay, L., Bondelid, T., Dewald, T., Johnston, J., Moore, R., Rea, A., 2012. NHDPlus Version
951 2: user guide. National Operational Hydrologic Remote Sensing Center, Washington, DC.
952
953 Mei, Y. and Anagnostou, E. N., 2015. A hydrograph separation method based on information
954 from rainfall and runoff records, J. Hydrol., 523, 636–649,
955 https://doi.org/10.1016/j.jhydrol.2015.01.083.
956
957 Merz, R., Blöschl, G., and Parajka, J., 2006. Spatio-temporal variability of event runoff
958 coefficients, J. Hydrol., 331, 591–604, https://doi.org/10.1016/j.jhydrol.2006.06.008.
959
960 Newman, A.J., Mizukami, N., Clark, M.P., Wood, A.W., Nijssen, B., Nearing, G., 2017.
961 Benchmarking of a physically based hydrologic model. J. Hydrometeorol. 18, 2215–2225,
962 http://doi.org/10.1175/JHM-D-16-0284.1.
963
964 Niu, G.-Y., Yang, Z.-L., Mitchell, K.E., Chen, F., Ek, M.B., Barlage, M., Kumar, A., Manning,
965 K., Niyogi, D., Rosero, E., Tewari, M., Xia, Y., 2011. The community Noah land surface
966 model with multiparameterization options (Noah-MP): 1. Model description and evaluation with
967 local-scale measurements. J. Geophys. Res. 116, D12109, doi:10.1029/2010JD015139.
968
969 NOAA National Weather Service, 2012. NWS Manual 10-950. Definitions and General
970 Terminology. Hydrological Services Program, NWSPD 10-9,
971 http://www.nws.noaa.gov/directives/sym/pd01009050curr.pdf.
972
973 R Core Team, 2019. R: A language and environment for statistical computing. R Foundation for
974 Statistical Computing, Vienna, Austria, https://www.R-project.org/.
975
976 Rathinasamy, M., Khosa, R., Adamowski, J., Ch, S., Partheepan, G., Anand, J., & Narsimlu, B.,
977 2014. Wavelet-based multiscale performance analysis: An approach to assess and improve
978 hydrological models. *Water Resources Research*, *50*(12), 9721–9737.
979
980 Schaefli, B., Zehe, E., 2009. Hydrological model performance and parameter estimation in the
981 wavelet-domain. Hydrol. Earth Syst. Sci. 13, 1921-1936, https://doi.org/10.5194/hess-13-1921-
982 2009.
983

984     Seibert, S. P., Ehret, U., Zehe, E. 2016. Disentangling timing and amplitude errors in streamflow
985     simulations. Hydrology and Earth System Sciences, 20, 3745–3763. https://doi.org/10.5194/hess-
986     20-3745-2016
987
988     Torrence, C., Compo, G.P., 1998. A Practical Guide to Wavelet Analysis. Bull. Am. Meteorol.
989     Soc. 79(1), 61-78.
990
991     Veleda, D., Montagne, R., Araujo, M., 2012. Cross-wavelet bias corrected by normalizing
992     scales. J. Atmos. Ocean. Technol. 29, 1401-1408.
993
994     Weedon, G.P., Prudhomme, C., Crooks, S., Ellis, R.J., Folwell, S.S., Best, M.J., 2015.
995     Evaluating the performance of hydrological models via cross-spectral analysis: case study of the
996     Thames Basin, United Kingdom. J. Hydrometeorol. 16(1), 214–231. http://doi.org/10.1175/JHM-
997     D-14-0021.1.