



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

3
4 Erin Towler^{1*} and James L. McCreight^{1,2}

5 ¹ National Center for Atmospheric Research (NCAR), P.O. Box 3000, Boulder, CO 80307

6 * Corresponding author: towler@ucar.edu, <https://orcid.org/0000-0002-1784-1346>

7 ² orcidid: 0000-0001-6018-425X

8 **Abstract**

9
10 Streamflow timing errors (in the units of time) are rarely explicitly evaluated, but are
11 useful for model evaluation and development. Wavelet-based approaches have been shown to
12 reliably quantify timing errors in streamflow simulations, but have not been applied in a
13 systematic way that is suitable for model evaluation. This paper provides a step-by-step
14 methodology that objectively identifies events, and then estimates timing errors for those events,
15 in a way that can be applied to large-sample, high-resolution predictions. Step 1 applies the
16 wavelet transform to the observations, and uses statistical significance to identify observed
17 events. Step 2 utilizes the cross-wavelet transform to calculate the timing errors for the events
18 identified in Step 1. The approach also includes a quantification of the confidence in the timing
19 error estimates. The methodology is illustrated using real and simulated stream discharge data
20 from several locations to highlight key method features. The method groups event timing errors
21 by dominant timescales, which can be used to identify the potential processes contributing to the
22 timing errors and the associated model development needs. For instance, timing errors that are
23 associated with the diurnal melt cycle are identified. The method is also useful for documenting
24 and evaluating model performance in terms of defined standards. This is illustrated by showing
25 version-over-version performance of the National Water Model (NWM) in terms of timing
26 errors.

27
28 **1. Introduction**

29



30 Common verification metrics used to evaluate streamflow simulations are typically
31 aggregated measures of model performance, e.g., the Nash Sutcliffe Efficiency (NSE) and the
32 related root mean square error (RMSE). Although typically used to assess errors in amplitude,
33 these statistical metrics include contributions from errors in both amplitude and timing (Ehret
34 and Zehe 2011), making them difficult to use for diagnostic model evaluation (Gupta et al.
35 2008). Furthermore, common verification metrics are calculated using the entire time series,
36 whereas timing errors require comparing localized features or events in the data. This paper
37 focuses explicitly on event timing error estimation, which is not routinely evaluated, despite its
38 potential benefit for model diagnostics (Gupta et al. 2008) and practical forecast guidance (Liu et
39 al. 2011).

40 The fundamental challenge with evaluating timing errors is identifying what constitutes
41 as an “event” in the two time series being compared. Identifying events is typically subjective,
42 time consuming, and not practical for large-sample hydrological applications (Gupta et al. 2014).
43 Most methods for identifying events have focused on flooding events. One common approach to
44 identifying flooding events is to use peak-over-threshold methods. The thresholds used for such
45 analyses are often either based on historical percentiles (e.g., the 95th percentile) or on local
46 impact levels (river stage), such as the National Weather Service (NWS) flood
47 categories (NOAA National Weather Service, 2012). Timing error metrics are often calculated
48 from the peaks of these identified events. For example, the Peak Time Error, or its derivative the
49 Mean Absolute Peak Time Error, requires matching observed and simulated event peaks, and
50 calculating their offset (Ehret and Zehe 2011). While this may be straightforward visually, it can
51 be difficult to automate; some of the reasons for this are discussed below.



52 Difficulties arise using thresholds for event identification. For example, exceedances can
53 cluster if a hydrograph vacillates above and below a threshold, begging the question: Is it one or
54 multiple events? Which peak should be used for the assessment? In the statistics of extremes,
55 declustering approaches can be applied to extract independent peaks (e.g., Coles 2001), but this
56 reductionist approach may miss relevant features. For instance, if background flows are elevated
57 for a longer period of time before and after the occurrence of these “events”, the threshold-based
58 analysis identifies features of the flow separately from the primary hydrologic process
59 responsible for the event. If one focuses just on peak timing differences in this example, that
60 timing error may only apply to some small fraction of the total flow of the larger event which
61 happens mainly below the threshold. Further, for overall model diagnosis that focuses on model
62 performance for all events, not just flood events, variable thresholds would be needed to account
63 for different kinds of events (e.g., a daily melt event versus a convective precipitation event).

64 Using a threshold-approach to identify events and timing error assessment, Ehret and
65 Zehe (2011) develop an intuitive assessment of hydrograph similarity, the Series Distance. This
66 algorithm is later improved upon by Siebert et al (2016). The procedure matches observed and
67 simulated segments (rise or recession) of an event, and then calculates the amplitude and timing
68 errors, as well as the frequency of event agreement. The Series Distance requires smoothing the
69 time series, identifying an event threshold, and selecting a time range to consider two segments
70 matching.

71 Liu et al (2011) developed a wavelet-based method for estimating model timing errors.
72 Although wavelets have been applied in many hydrologic applications such as model analysis
73 (e.g. Lane 2007; Weedon et al. 2015; Schaeffli and Zehe 2009, Rathinasamy et al. 2014) and
74 post-processing (Bogner and Kalas 2007; Bogner and Pappenberger 2011), Liu et al. were the



75 first to use it for timing error estimation. Liu et al. (2011) apply a cross-wavelet transform
76 technique to streamflow time series for 11 headwater basins in Texas. Timing errors are
77 estimated for medium- to-high flow “events” that are determined a priori by threshold
78 exceedance. They use synthetic as well as real streamflow simulations to test the utility of the
79 approach. They show that the technique can reliably estimate timing errors, though they
80 conclude that it is less reliable for multi-peak or consecutive “events” (defined qualitatively).
81 ElSaadani and Krajewski (2017) followed the cross-wavelet approach used by Liu et al (2011) to
82 provide similar analysis and further investigate the effect of the choice of mother wavelet on the
83 timing error analysis. Ultimately, they recommended that in the situation of multiple, adjoining
84 flow peaks the improved time localization of the Paul wavelet might justify its poorer frequency
85 localization compared the Morlet wavelet.

86 Liu et al. (2011) provide a starting point for the work in this paper where we develop two
87 new bases for their method: 1) objective event identification for timing error evaluation and 2)
88 the use of observed events as the basis for the model timing error calculations The latter is
89 important for “model benchmarking”, i.e., the practice of evaluating models in terms of defined
90 standards (e.g., Luo, et al. 2012; Newman et al. 2017). Here, the use of observed events provides
91 a baseline by which to evaluate changes and to compare multiple versions or experimental
92 designs.

93 This paper provides a methodology for using wavelet analysis to quantify timing errors in
94 hydrologic simulations. Our contribution is a systematic approach that integrates 1) statistical
95 significance to identify events with 2) a basis for timing error calculations independent of model
96 simulations (i.e., benchmarking). We apply our method to evaluation of high-resolution
97 streamflow prediction. The paper is organized as follows: Section 2 provides an overview of the



98 conceptual approach of using wavelets to identify events and estimate timing errors, and Section
99 3 provides the detailed methodology. In Section 4, we describe the software and data, as well as
100 provide a simple illustration of the method using real and simulated streamflow data. In Section
101 5, we provide results, including select examples to highlight features of the method and version-
102 over-version comparisons. Section 6 is the discussion and conclusions, including how specific
103 choices may vary by application.

104 **2. Conceptual Overview**

105 Before going into technical details of the Method (Section 3), we provide a conceptual
106 overview of the approach of using wavelets to identify events and estimate timing errors. We
107 provide a nomenclature table (Supplemental Table 1) of key terms relevant to the approach. The
108 wavelet transform (WT) expands the dimensionality of the original time series by introducing the
109 timescale (or period) dimension and returns power as a function of both time and timescale (e.g.
110 Torrence and Compo, 1998). This is illustrated in Figure 1: the streamflow time series (panel a) is
111 expanded into a 2-dimensional wavelet power spectrum (panel b). Where traditional model errors,
112 such as the aforementioned RMSE or NSE, reduce the information of the time series to a single
113 statistic, wavelet analysis expands the input signal and provides information on the dominant
114 timescales of the time series at each time. Wavelet analysis can therefore detect localized signals
115 in time series (Daubechies 1990), including hydrologic time series, which are often irregular or
116 aperiodic (i.e., events may be isolated and don't regularly repeat) or non-stationary. We note that
117 in many wavelet applications, timescale is referred to as "period". To emphasize that our study is
118 more focused on irregular events and less on periodic behavior of time series, we use the term
119 "timescale". The wavelet transform is the foundation of the view in this paper that events have



120 characteristics of both time and timescale. Timing errors, calculated from events defined this way,
121 therefore have dimensions of both time and timescale as well.

122 In their seminal wavelet study, Torrence and Compo (1998) outline a method for
123 objectively identifying statistical significance in the wavelet transform. We adopt this approach
124 and define “events” in the observed time series via statistical significance of the wavelet power
125 spectrum. The details are provided in the next section, however Figure 1 illustrates that the
126 events in the input time series (panel a) are defined as regions of the wavelet power spectrum
127 shown in panel b: events are inside the black contours ($\geq 95\%$ confidence level) but not inside
128 the cone of influence (regions where the colors are muted, this is explained in detail in Section
129 3). The wavelet power spectrum is only shown for the events in panel c. Events defined in this
130 way are a function of both time and timescale. Note that at a given time, events of different
131 timescales can occur simultaneously. What one may subjectively interpret as a single event in the
132 input time series is generally quantified by this definition as multiple coincident events at a
133 variety of timescales each with a different power (e.g. Figure 1, panel c). Although for some
134 locations there may be physical reasons to expect certain timescales to be important (e.g.,
135 seasonal cycle of snowmelt), the most important scales at which hydrologic signals occur at a
136 particular location are not necessarily known a priori. The wavelet power can be examined
137 across events to identify the most dominant, or what we call “characteristic” timescales for a
138 given time series; the procedure for this is detailed later in the technical methodological section
139 (Section 3.1.3). This approach to event detection is objective, data-driven, and portable across
140 diverse locations, which is important for large-sample hydrologic applications. We point out that
141 in the objective identification of events, we are not limited to flooding events. Rather, events are



142 defined more broadly: an event is when the wavelet power falls outside its standard statistical
143 power. This can be further subset into flooding events if desired.

144 Once observed events are identified by the method, we can calculate timing errors
145 between observed and simulated time series. The cross-wavelet timing error approach of Liu et
146 al (2011) is used, but we restrict our calculation of timing errors to the aforementioned regions of
147 statistically significant wavelet power in the observations; i.e., we calculate timing errors in
148 terms of *observed* events (Figure 1c). Because both the phase (timing error) and the significance
149 of the cross wavelet transform (XWT) computed between the observed and modeled time series
150 depends on the modeled time series, we use the observed event definition (Figure 1c) in the
151 calculation of the timing errors to provide a common, consistent basis independent of the models
152 evaluated (i.e., benchmarking). The portions of the observed wavelet spectrum used for
153 comparison may further be restricted depending on the analysis goals.

154 **3. Method for evaluating event timing errors**

155 This section provides the technical description of the methodology, and the steps can be
156 seen in an accompanying flowchart (Supplemental Figure 1).

157 *3.1. Step 1. Identify observed events*

158 The first step towards evaluating timing errors is to identify a set of observed events for
159 which the timing error should be calculated. We break this step into three sub-steps: 1a. Apply
160 the wavelet transform to observations, 1b. Determine all observed events using significance
161 testing, and 1c. Sample observed events to an event-set relevant to analysis.

162 *3.1.1. Step 1a. Apply wavelet transform to observations*

163 First, we apply the continuous wavelet transform to the observed time series. We provide
164 an overview of the main steps and equations for the wavelet transform here, though the reader is
165 referred to Torrence and Compo (1998) and Liu et al. (2011) for more details.



166 Before applying the WT, a mother wavelet needs to be selected. In Torrence and Compo
167 (1998), they discuss the key factors that should be considered when choosing the mother
168 wavelet. There are four main considerations, including (i) orthogonal or nonorthogonal, (ii)
169 complex or real, (iii) width, and (iv) shape. In this study, we follow Liu et al. (2011) in selecting
170 the nonorthogonal and complex Morlet wavelet:

171
$$\psi(n) = \pi^{-1/4} e^{iw_0 n} e^{-n^2/2},$$

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

173 Once the mother wavelet is selected, the WT is applied to a time series x_n , where n goes
174 from $n=0$ to $n=N-1$, with a time step of δt . The WT is the convolution of the time series with the
175 mother wavelet that has been scaled and normalized:

176
$$W_n(s) = \sum_{n'=0}^{N-1} x_{n'} \psi^* \left[\frac{(n'-n)\delta t}{s} \right],$$

177 where s is the scale parameter, the asterix indicates the complex conjugate of the wavelet
178 function. The wavelet power is defined as $|W_n^2|$. We use the bias corrected wavelet power
179 spectrum (Liu et al. 2007; Veleda et al. 2012), which ensures spectral peaks are comparable
180 across timescales. We also identify a maximum timescale that corresponds to our application.
181 We select 256 hours (~10 days), but this number could be higher or lower for other applications
182 and there are no real penalties for using too high a maximum (lower than the annual cycle).

183 Because we are applying the WT to a finite time series, there are timescale-dependent
184 errors at the beginning and end times of the power spectrum. These are referred to as the cone of
185 influence or COI (Torrence and Compo, 1998). We ignore all results within the COI in this
186 study.

187 3.1.2. Step 1b. Determine all observed events using significant testing

188 Once the WT is applied, the 2-dimensional (2-D) wavelet power spectra shows how the
189 features of the time series vary with both time and timescale. To identify areas of significance,



190 we apply Torrence and Compo's (1998) approach that compares the WT power spectra with a
191 power spectra from a red noise process. Specifically, the observed time series is fitted with an
192 order 1 autoregressive (AR1, or red noise) model, and the WT is applied to the AR1 time series.
193 The power spectra of the AR1 model provide the basis for the statistical significance testing.
194 Significance is determined if the power spectra are statistically different using a chi-squared test
195 with 95% confidence.

196 Statistical significance indicates an "event" at a given time and timescale: that is, the
197 wavelet power falls outside its standard statistical power. The result is the set of all events, i.e.,
198 each event is a combination of time and timescale (i.e., locations on the 2-D grid). We refer to
199 contiguous regions of statistical significance (in time and timescale) as "event clusters" (note that
200 no statistical clustering is performed).

201 3.1.3. Step 1c. Sample observed events to an event-set relevant to analysis

202 Step 1b results in the identification of all events at all timescales and times. In this sub-
203 step, the event space is sampled to suit the particular evaluation. Because the goal of this paper is
204 to evaluate model timing errors over long simulation periods, we choose to sample the event
205 space based on dominant timescales in the time-averaged observed wavelet spectra. For our
206 application we choose to further sub-sample the observed wavelet spectra by selecting, for each
207 characteristic timescale, the most powerful event within each event cluster. This is articulated in
208 the following bullets:

- 209 • *Calculate the average event power across each timescale:* Considering only the
210 statistically significant areas of the observed wavelet spectrum, calculate the average
211 power across each timescale.

212



- 213 • *Identify timescales of absolute and local average power maxima:* By plotting the average
214 power versus the timescale, the local and absolute maximums for average power can be
215 determined. The timescales corresponding to the absolute and local maxima of the
216 average power are called the characteristic timescales of the observed wavelet spectrum.
217 This is the first subset of events: all events that fall within the characteristic time scales.
218
219 • *Identify events with maximum power for each event cluster:* As previously mentioned,
220 events can also be grouped into “event clusters”, that is, contiguous significant areas. We
221 can use this to further sample from the event-set created in the last bullet: across each
222 characteristic timescale, we identify the event with maximum power for each event
223 cluster. This is the second event subset: all events with maximum power for each cluster
224 that fall within a characteristic timescale.

225
226 *3.2. Step 2. Calculate Timing Errors*

227 Step 1 identifies events by applying a wavelet transform to the observed time series. To
228 calculate the timing error of a modeled time series, we perform its cross wavelet transform with
229 the observed time series, as detailed in this section.

230 *3.2.1. Step 2a. Apply cross-wavelet transform (XWT) to observations and simulations*

231 Given the WT of an observed time series $W_n^X(s)$ and a modeled time series $W_n^Y(s)$, the
232 cross-wavelet spectrum can be defined as:

233
$$W_n^{XY}(s) = W_n^X(s)W_n^{Y*}(s),$$

234 where the asterix implies the complex conjugate. The cross-wavelet power is defined as
235 $|W_n^{XY}(s)|$.

236 Similar to Step 1b of the WT, we can also calculate the areas of significance for the
237 XWT. These are not the same as the areas of significance for the WT. The significant areas of



238 the XWT vary with each simulation, and are therefore not useful for evaluation on their own.
239 Nevertheless, we are interested in the overlap between the significant areas of the observed WT
240 and the significant areas of the cross-wavelet transform, and this is used to quantify our
241 confidence in the timing error estimate. We discuss this further in Step 2d.

242 3.2.2. Step 2b. Calculate the cross-wavelet timing errors

243 To calculate the timing errors, we first compute the phase angle of the cross-wavelet
244 spectrum. The phase angle gives the phase difference and can be computed as:

$$245 \phi_n^{XY}(s) = \tan^{-1} \left[\frac{\Im((s^{-1}W_n^{XY}(s)))}{\Re((s^{-1}W_n^{XY}(s)))} \right],$$

246 where \Im is the imaginary and \Re is the real component of $W_n^{XY}(s)$.

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

$$248 \Delta t_n^{XY}(s) = \phi_n^{XY}(s) * T/2\pi,$$

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

250 3.2.3. Step 2c. Subset cross-wavelet timing errors to sampled observed events

251 Step 2b results in an estimate of timing errors for all times and timescales in the cross-
252 wavelet transform space. In our application we are interested in the timing errors that correspond
253 to the identified sample of *observed* events, especially for events at the characteristic timescales
254 (the first event-set in step 1c) and for the maximum power events in each cluster (the second
255 event-set in step 1c). The latter provides a single timing error for each event cluster at each
256 characteristic timescale, which could be used in a post-processing step to provide a cluster-by-
257 cluster timing correction, if desired.

258 It is important to point out that for other applications, there could be other ways to
259 interrogate the timing errors that result from the cross-wavelet transform. Some of these
260 possibilities are noted in the Discussion section.

261 3.2.4. Step 2d. Quantify the confidence in the timing error estimate



262 To interpret our confidence in the timing error estimate, we can examine the overlap
263 between the significant areas of the observed WT and the significant areas of the XWT.

264 We can look at percent (%) overlap, that is, how many of the XWT events overlap with
265 the WT events, either for all events or for the sampled event-sets. An overlap close to 0% would
266 indicate that the model did not do a good job of simulating the observations – or it is a “miss”
267 (flood is observed but not forecasted). If the overlap was 100%, it would be close to a perfect
268 simulation. Second, if we are looking at a single timing error for each event cluster, we may look
269 to see if that event is significant in the XWT. If it is not, it gives us less confidence in the
270 estimate.

271 We note that because we are calculating timing errors in terms of observed events, there
272 is no information about “false alarms”, where a flood is forecasted but not observed.

273 **4. Application of the Framework**

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

276 *4.1. Data*

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

284 NOAA’s National Water Model (NWM,
285 <https://www.nco.ncep.noaa.gov/pmb/products/nwm/>) is an operational model that produces



286 hydrologic analyses and forecasts over the continental United States (CONUS) and Hawaii (as of
287 version 2.0). The model is forced by downscaled atmospheric states and fluxes from NOAA's
288 operational weather models. Next, the NoahMP (Niu et al 2011) land surface model calculates
289 energy and water states and fluxes. Water fluxes propagate down the model chain through
290 overland and subsurface (soil and aquifer representations) water routing schemes to reach a
291 stream channel model. The NWM applies the three parameter Muskingum-Cunge river routing
292 scheme to a modified version of the NHD-Plus version 2 (McKay et al. 2012) river network
293 representation.

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

300 To apply the methodology, we note that the observed and simulated datasets must be
301 paired (overlapping). Further, for evaluation, any new simulation must also be paired with the
302 observed. Missing data, which is common in observed time series, can be problematic and can
303 result in false significance. We account for this our methodology by calculating the XT and
304 XWT on each complete time series. This will be illustrated in the forthcoming example at Taylor
305 River, CO.

306 *4.2. Application*

307 For illustration purposes we apply Steps 1 and 2 to an observed time series in Onion
308 Creek, TX; for simplicity, we select an isolated peak (Figure 1a). First, we apply the wavelet
309 transform to the observations (Figure 1b). This shows the time series in terms of its power by



310 time and timescale, with warmer colors indicating more power. The black outline shows the
311 areas of significance and the muted colors indicate the COI. To determine all observed events,
312 we identify all the points that are significant and outside the COI (Figure 1c). Next, we average
313 the power across each timescale: to the right of Figure 1b we show power averaged across all
314 points for each timescale, and to the right of Figure 1c we show power averaged across just the
315 events for each timescale. The latter is the one used to identify our characteristic scales. In this
316 case, there is a single maximum at 22 hours. For the characteristic timescale, we see there is only
317 1 event cluster and the event with maximum power is marked with a star (Figure 1d).

318 For Step 2, we use the same Onion Creek, TX, peak from Figure 1a, and add a prescribed
319 timing error of +5 hours to every point in the original time series (Figure 2a) to create a synthetic
320 time series. We perform the cross-wavelet transform between the observed and synthetic time
321 series (Figure 2b). The arrows in Figure 2b indicate the phase offset, which are used to calculate
322 the timing error (Figure 2c). The timing error estimates show that for timescales greater than 10
323 hours, we get back the prescribed timing error of 5 hours, i.e., the scale must be at least double
324 the timing error. In this case, because we are adding a prescribed error, the error is approximately
325 5 hours for all events, including for the characteristic timescale of 22 hours.

326 Finally, we repeat Step 2, but compare the observation of this event to actual model data
327 from NWM V1.2. This shows that the model is early (Supplemental Figure 2a). We perform the
328 cross wavelet transform (Supplemental Figure 2b) and examine the timing error (Supplemental
329 Figure 2c). Table 1 summarizes the results: the mean error across the 22-hour characteristic
330 timescale is -3.2 hours, as is the error for the cluster's maximum power. All events in the cluster
331 are also significant in the XWT (100%), and the cluster maximum is also significant, providing
332 confidence in this timing error estimation.



333

334 **5. Results**

335 In this section, modeled data is used from several locations and time series to highlight the
336 features of the method, finishing with version-over-version comparisons to illustrate the utility
337 for evaluation.

338 5.1. Pemigewasset River, NH

339 This example uses time series from the Pemigewasset River, NH. First, we examine a three-
340 month time series that exhibits multiple peaks above a base flow (Figure 3a). By eye, it is fairly
341 straightforward to pick out three main peaks. The wavelet transform (Figure 3b and 3c) reveals
342 up to three event clusters, depending on the characteristic timescale examined (Figure 3d). When
343 we plot the average power by timescale (right of Figure 3c), we see that there are nine relative
344 maxima (small grey dots) – hence there are 9 characteristic scales for this example.

345 In Step 2, we compare the same time series with output from NWM V1.2 (Figure 4a), apply
346 the cross-wavelet transform (Figure 4b), and calculate the timing error for all observed events
347 (Figure 4c). As previously mentioned, we are interested in the timing errors corresponding to
348 observed events at the characteristic timescales. In Figure 5a, the panels are ordered by
349 timescales from highest to lowest average power; we only show the top 5 characteristic scales,
350 using the first-subset of events, grouped by cluster. The first panel, where timescale = 24.8 hours,
351 is the absolute maximum. This shows two cluster distributions: for cluster one, the model is late
352 for most events, and cluster two shows the model is early; the dark shading indicates that most of
353 the events are significant in the XWT. The next two dominant scales have similar average power
354 and are of the same order of magnitude at 27.8 hours and 33.1 hours; if we had applied
355 smoothing to the graph of average power by timescale, these relative maxima would smooth out.
356 We will revisit this in the Discussion, when we discuss pathways to implementation. The



357 characteristic scale with the next highest maxima occurs at 111 hours, which is a different order
358 of magnitude, suggesting that this may have a different physical process driving it. This shows
359 the model to be late for both clusters, and results are similar for a timescale of 148 hours. We
360 don't show results for the remaining 4 characteristic time scales with lower average power, since
361 they have similar characteristic timescale values and associated timing errors to what has already
362 been shown.

363 We can see how looking at the timing errors using the cluster distributions will get harder as
364 the number of clusters increase, so it is also useful to summarize the information by looking at
365 each cluster mean and max. If we run the methodology on the full 5-year Pemigewasset River
366 time series, we can compare the mean and max timing errors for each characteristic time scale
367 using box plots where the outline is shaded by the average confidence (Supplemental Figure 3).
368 Table 2 summarizes this information. For example, the absolute maxima, at the 17.5 hour
369 timescale has 86 clusters, and a timing error centered around zero (-0.43 hours), 75% of which
370 are significant in the XWT. This is very similar to the results for the cluster max, as it is for the
371 rest of the characteristic time scales. One other thing to note is that as expected, because the
372 characteristic time scales are data driven, they are not the same as they were for the 3-month
373 period.

374 5.2. Bad River, SD

375 The second example uses a two-month time series from the Bad River, SD, to illustrate the
376 concept of consecutive peaks (Figure 6a). Whereas in the previous example it was fairly
377 straightforward to pick out 3 distinct peaks, in this time series, there is one noticeable peak
378 centered around June the 1st, with smaller peaks preceding and following it. The question is
379 whether or not this is one event cluster or multiple? Looking at the wavelet transform (Figure 6b



380 and 6c), we can see that for smaller timescales, there are more clusters, but for longer timescales,
381 they are considered a single cluster.

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

388 This example has another interesting feature: namely that there is a false alarm in the model
389 just before July 15. We note that because of our methodology, there is no observed event at that
390 time, and therefore no timing error to be calculated, that is there is no information in the timing
391 error statistics in terms of false alarms.

392 5.3. Taylor River, CO

393 In this example, we will examine a time series from Taylor River, CO, that illustrates peaks
394 that are driven by different processes. The Taylor River is in a mountainous area where the
395 spring hydrology is dominated by snowmelt runoff. To start, we will look at a portion of the
396 spring melt season, where we can visibly see a diurnal signal (Figure 8). However, while it's
397 easy to see that the model is too high in amplitude, it's hard to visually tell much about the
398 timing error. Figure 9 shows that for the characteristic time scale of 23.4 hours, the model is
399 usually early, with high confidence.

400 Supplemental Figure 4a shows a year-long time series from Taylor River, CO, where we can
401 see the snowmelt runoff in spring, but also several peaks in summer, likely driven by summer
402 rains. From the WT, we again see the peak in the characteristic time scales at about 24 hours
403 (right of Supplemental Figure 4c), but there is another maxima at 99 and 118 hour timescales,



404 relating to flows from the summer rains. Looking at Figure 10, starting with the 24 hour
405 timescale, we see that for the clusters that are significant in the XWT, the model is generally
406 early. For the 118 and 99 hour timescale, the model is also early, but those cluster events are not
407 statistically significant in the XWT. This suggests that we are confident in the early timing errors
408 of the model for the diurnal snowmelt cycle, and this could be used as qualitative guidance for
409 model performance at this site until the model performance is improved. However, we show that
410 it is less reliable for the early timing errors for the summer peaks. This underscores the key point
411 that timing errors are timescale dependent, and can help diagnose which processes to target for
412 improvements.

413 Supplemental Figure 4b also illustrates how missing data is handled: this results in additional
414 COIs (muted colors) to account for the edge effects, and areas of the COI are ignored in our
415 analyses.

416 5.4. Evaluating Model Performance

417
418 Finally, we show how the methodology can be used for evaluating performance changes
419 across NWM versions. We point out that none of the NWM version upgrades were targeting
420 timing errors, so these results just provide a demonstration. We use a 5-year overlapping time
421 series and cluster max for the results, but cluster mean results were similar (not shown). For the
422 NWM V1.0 for Onion Creek, we see that for the 29.5 hour timescale, there were 17 clusters, for
423 which the median timing error is -1.4 hours, and all were significant in the XWT (Table 3).
424 Comparing V1.0, V1.1, and V1.2, the results for Onion Creek show that the median timing error
425 has gotten slightly earlier (worse), although the distribution became tighter from V1.0 to V1.1
426 and V1.2 (Figure 11). In Figure 11, the dark blue color of the boxplot outline indicates that there
427 is high confidence in the timing error, as the overlapping significance is close to 100% for the



428 top three characteristic timescales. Using the 5-year overlapping time series for Pemigewasset
429 River, NH, we see that the median timing error improved by getting closer to zero, but that the
430 distribution became wider (Figure 12). Again, the confidence is fairly high (>80%) across
431 characteristic time scales and versions (Table 4), and >60 clusters were used in the estimations.
432 Using 5-years from Taylor River, CO (Supplemental Table 2, Supplemental Figure 5), we see
433 that for the characteristic scale of 235 hours (~10 days), has low confidence (~25%) for the 4
434 sampled clusters; the timescale of 23.4 hours has a median that is consistently early by around 6
435 hours, with the version model confidence ranging from 44% to 67% (Supplemental Table 2).
436 Results for the Bad River can be seen in Supplemental Table 3 and Supplemental Figure 6.

437
438
439

6. Discussion and Conclusions

440 In this paper, we develop a systematic, data-driven methodology to objectively identify
441 events and estimate timing errors in large-sample, high-resolution hydrologic models. The
442 method was developed towards several intended uses: Primarily, it was developed for model
443 evaluation, so that model performance can be documented in terms of defined standards. We
444 illustrate this with the version-over-version NWM comparisons. Second, it can be used for model
445 development, whereby potential timing error sources can be diagnosed and targeted for
446 improvement. Related to this point, given the advantages of calibrating using multiple-criteria
447 (e.g., Gupta et al. 1998), timing errors could be used as part of a larger calibration strategy.
448 However, as noted in the consecutive peaks example for the Bad River, minimizing timing errors
449 at one timescale may not translate to improvements in timing errors (or other metrics) at other
450 scales. Wavelet analysis has also been used directly as an objective function for calibration,
451 although a difficulty is in determining the similarity measure to use (e.g. Schaeffli and Zehe 2009,
452 Rathinasamy et al. 2014). Future research will investigate the properties of timing errors for



453 calibration. Finally, the approach can be used for model interpretation, as estimating timing
454 errors provides a characterization of the uncertainty (i.e., for a given timescale, the model is
455 generally late or early), as well as a measure of the confidence, that could be useful for
456 qualitative forecast guidance.

457 Given the fact that several subjective choices were made specific to our application and
458 goals, we think it is important to highlight that we have made the analysis framework openly
459 available (detailed in the code availability section below), so the method can be extended or
460 refined by the community right away. For instance, because of our focus on model evaluation
461 and development, we use the observed WT to identify events. However, in other instances it
462 might be sufficient to only sample events that are in the significant areas of the XWT (essentially
463 to identify the characteristic scales and event-set directly from the XWT instead of from the
464 WT). This might be reasonable for applications that are more focused on model interpretation in
465 a real-time forecasting mode, but it would not allow for version comparison and it is not
466 guaranteed that all the important characteristic scales would be identified (i.e., the model may
467 not capture some real-world processes, and therefore miss the associated characteristic
468 timescales). We only look at the timing errors from an event-set relevant to our analysis, but
469 there are other ways to subset the events that might be more suitable to other applications. For
470 instance, we define the event set broadly, but it could be subset for high peak or flooding events
471 to compare with traditional peak-over-threshold approaches. For example, Supplemental Figure
472 7 shows the maximum streamflows for the event-set from the 5 year run at Taylor River; this
473 event-set could be filtered to include only events above a given threshold. The method provides a
474 quantification of the confidence in the timing errors, and we include all timing errors in our
475 summaries. However, it might make more sense to drop points that do not have a high



476 confidence (i.e., with a low percent of events that significantly overlap between the XT and the
477 XWT) and to flag those events as misses.

478 Another point that arises is how many characteristic timescales should be examined.
479 Here, we average the power across timescales and identify characteristic scales to be at every
480 absolute and relative maxima. As seen in the illustrative examples, this can result in multiple
481 characteristic scales, some of which can be quite similar, suggesting that events at those scales
482 are from similar or related processes. One solution could be to smooth the average power by
483 timescale, which would reduce the number of local maxima, or to look at timing errors within a
484 band of timescales. It is also important to note that the characteristic scales are data-driven, so
485 they will change with different lengths of observed time series. Longer runs capture more events
486 and should converge on the more dominant timescales and events for a location. However, for
487 performance evaluation, overlapping time periods are needed.

488 In our application of the WT, we follow Liu et al. (2011) and select the Morlet as the
489 mother wavelet. However, results are sensitive to the mother wavelet selected. Further discussion
490 of mother wavelet choices can be found in Torrence and Compo (1998) and in ElSaadani and
491 Krajewski (2017).

492 In short, this paper provides a systematic, flexible, and computationally efficient
493 methodology that is appropriate for model evaluation and comparison, and is useful for model
494 development and guidance. Future work will apply the approach to identify characteristic
495 timescales across the United States, as well as to assess the associated timing errors in the NWM.

496 **Code/Data Availability**

497 The code for reproducing the figures in this paper as well as extended vignettes/notebooks are
498 provided in public github repository https://github.com/NCAR/wavelet_timing. In addition to



499 reproducing the analyses and figures in this paper, several jupyter notebooks provide more
500 detailed analyses of the time series included in this paper. We emphasize that the analysis
501 framework is meant to be flexible and adapted to similar applications where different statistics
502 may be desired. The figures created are specific to the applications in this paper but provide a
503 starting point for other work.

504 The core code is provided in the public “rwrflhydro” R package
505 <https://github.com/NCAR/rwrflhydro>. The package can be installed as described by the
506 README document in the repository and in the Supplemental Online Materials for this paper.
507 The code is written in the open-source R language (R Core Team 2019) and builds off multiple,
508 existing R packages. Most notably the wavelet and cross-wavelet analyses are performed using
509 the “biwavelet” package (Gouhier et al. 2018).

510 **Credit Author Statement**

511 ET and JLM collaborated to develop the methodology. ET led the results analysis and
512 manuscript preparation and revisions. JLM developed the initial idea for the work, the open
513 source software, and visualizations.

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

515

516 **Acknowledgements:**

517

518 The authors would like to thank Dave Gochis for useful discussions and Aubrey Dugger
519 for providing NWM data. We thank the NOAA/OWP and NCAR NWM team for its support of
520 this research. This research is funded by the NOAA Office of Water Prediction and the Joint
521 Technology Transfer Initiative grant 2018-0303-1556911. This material is based upon work
522 supported by the National Center for Atmospheric Research (NCAR), which is a major facility



523 sponsored by the National Science Foundation (NSF) under Cooperative Agreement No.

524 1852977.

525

526 **References:**

527

528 Bogner, K., Pappenberger, F., 2011. Multiscale error analysis, correction, and predictive
529 uncertainty estimation in a flood forecasting system. *Water Resour. Res.* 47, W07524,
530 doi:10.1029/2010WR009137.

531

532 Coles, S., 2001. *An Introduction to Statistical Modeling of Extreme Values*, Springer Ser. Stat.
533 Springer, London.

534

535 Daubechies, I., 1990. The wavelet transform time-frequency localization and signal analysis.
536 *IEEE Trans. Inform. Theory* 36, 961-1004.

537

538 ElSaadani, M., Krajewski, W. F., 2017. A time-based framework for evaluating hydrologic
539 routing methodologies using wavelet transform. *Journal of Water Resource and Protection*, 9(7),
540 723–744.

541

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

545

546 Gouhier, T.C., Grinsted, A., Simko, V., 2018. R package biwavelet: Conduct Univariate and
547 Bivariate Wavelet Analyses (Version 0.20.17), <https://github.com/tgouhier/biwavelet>.

548

549 Gupta, H.V., Kling, H., Yilmaz, K.K., Martinez, G. F., 2009. Decomposition of the mean
550 squared error and NSE performance criteria: Implications for improving hydrological modelling.
Journal of hydrology, 377(1-2), 80-91.

551

552 Gupta, H.V., Perrin, C., Blöschl, G., Montanari, A., Kumar, R., Clark, M., Andréassian, V.,
553 2014. Large-sample hydrology: a need to balance depth with breadth. *Hydrol. Earth Syst. Sci.*
554 18, 463-477, <https://doi.org/10.5194/hess-18-463-2014>.

555

556 Gupta, H.V., Sorooshian, S., Yapo, P.O., 1998. Towards improved calibration of hydrologic
557 models: multiple and non-commensurable measures of information. *Water Resources Research*
558 34(4): 751–763.

559

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

563

564 Lane, S.N., 2007. Assessment of rainfall–runoff models based upon wavelet analysis. *Hydrol.*
565 *Processes* 21, 586–607, <https://doi.org/10.1002/hyp.6249>.

566



- 567 Liu, Y., Liang, X.S., Weisberg, R.H., 2007. Rectification of the bias in the wavelet power
568 spectrum. *J. Atmos. Oceanic Technol.* 24, 2093–2102.
569
- 570 Liu, Y., Brown, J., Demargne, J., Seo, D. J., 2011. A wavelet-based approach to assessing timing
571 errors in hydrologic predictions. *J. Hydrol.* 397(3–4), 210–224.
572 <http://doi.org/10.1016/j.jhydrol.2010.11.040>
573
- 574 Luo, Y. Q., Randerson, J.T., Abramowitz, G., Bacour, C., Blyth, E., Carvalhais, N., ... Orme, C.
575 E., 2012. A framework for benchmarking land models. *Biogeosciences* 9, 3857–3874,
576 <https://doi.org/10.5194/bg-9-3857-2012>.
577
- 578 McKay, L., Bondelid, T., Dewald, T., Johnston, J., Moore, R., Rea, A., 2012. NHDPlus Version
579 2: user guide. National Operational Hydrologic Remote Sensing Center, Washington, DC.
580
- 581 Newman, A.J., Mizukami, N., Clark, M.P., Wood, A.W., Nijssen, B., Nearing, G., 2017.
582 Benchmarking of a physically based hydrologic model. *J. Hydrometeorol.* 18, 2215–2225,
583 <http://doi.org/10.1175/JHM-D-16-0284.1>.
584
- 585 Niu, G.-Y., Yang, Z.-L., Mitchell, K.E., Chen, F., Ek, M.B., Barlage, M., Kumar, A., Manning,
586 K., Niyogi, D., Rosero, E., Tewari, M., Xia, Y., 2011. The community Noah land surface
587 model with multiparameterization options (Noah-MP): 1. Model description and evaluation with
588 local-scale measurements. *J. Geophys. Res.* 116, D12109, doi:10.1029/2010JD015139.
589
- 590 NOAA National Weather Service, 2012. NWS Manual 10-950. Definitions and General
591 Terminology. Hydrological Services Program, NWSPD 10-9,
592 <http://www.nws.noaa.gov/directives/sym/pd01009050curr.pdf>.
593
- 594 R Core Team, 2019. R: A language and environment for statistical computing. R Foundation for
595 Statistical Computing, Vienna, Austria, <https://www.R-project.org/>.
596
- 597 Rathinasamy, M., Khosa, R., Adamowski, J., Ch, S., Partheepan, G., Anand, J., & Narsimlu, B.,
598 2014. Wavelet-based multiscale performance analysis: An approach to assess and improve
599 hydrological models. *Water Resources Research*, 50(12), 9721–9737.
600
- 601 Schaeffli, B., Zehe, E., 2009. Hydrological model performance and parameter estimation in the
602 wavelet-domain. *Hydrol. Earth Syst. Sci.* 13, 1921–1936, [https://doi.org/10.5194/hess-13-1921-](https://doi.org/10.5194/hess-13-1921-2009)
603 2009.
604
- 605 Seibert, S. P., Ehret, U., Zehe, E. 2016. Disentangling timing and amplitude errors in streamflow
606 simulations. *Hydrology and Earth System Sciences*, 20, 3745–3763. [https://doi.org/10.5194/hess-](https://doi.org/10.5194/hess-20-3745-2016)
607 [20-3745-2016](https://doi.org/10.5194/hess-20-3745-2016)
608
- 609 Torrence, C., Compo, G.P., 1998. A Practical Guide to Wavelet Analysis. *Bull. Am. Meteorol.*
610 *Soc.* 79(1), 61–78.
611
- 612 Veleda, D., Montagne, R., Araujo, M., 2012. Cross-wavelet bias corrected by normalizing

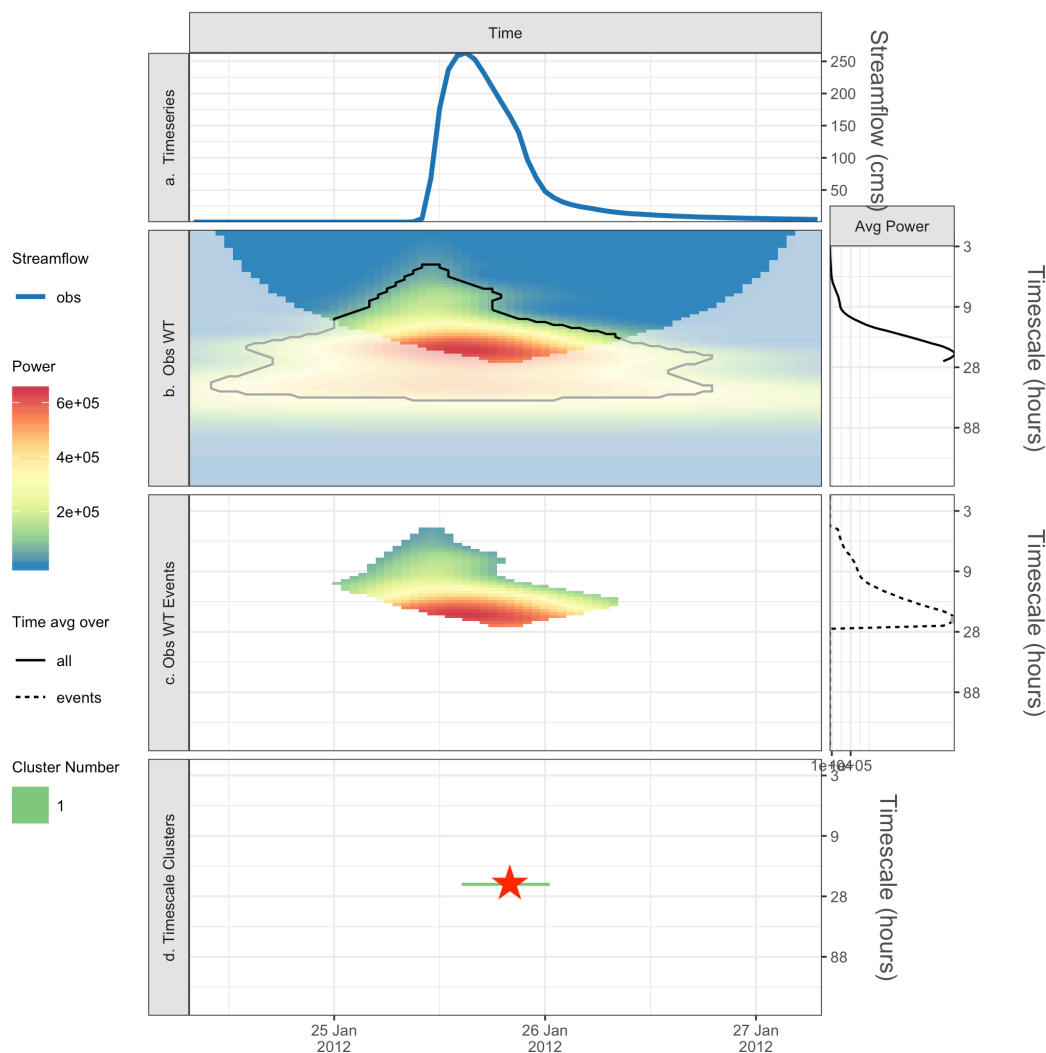


613 scales. J. Atmos. Ocean. Technol. 29, 1401-1408.
614
615 Weedon, G.P., Prudhomme, C., Crooks, S., Ellis, R.J., Folwell, S.S., Best, M.J., 2015.
616 Evaluating the performance of hydrological models via cross-spectral analysis: case study of the
617 Thames Basin, United Kingdom. J. Hydrometeorol. 16(1), 214–231. [http://doi.org/10.1175/JHM-](http://doi.org/10.1175/JHM-D-14-0021.1)
618 [D-14-0021.1](http://doi.org/10.1175/JHM-D-14-0021.1).
619
620

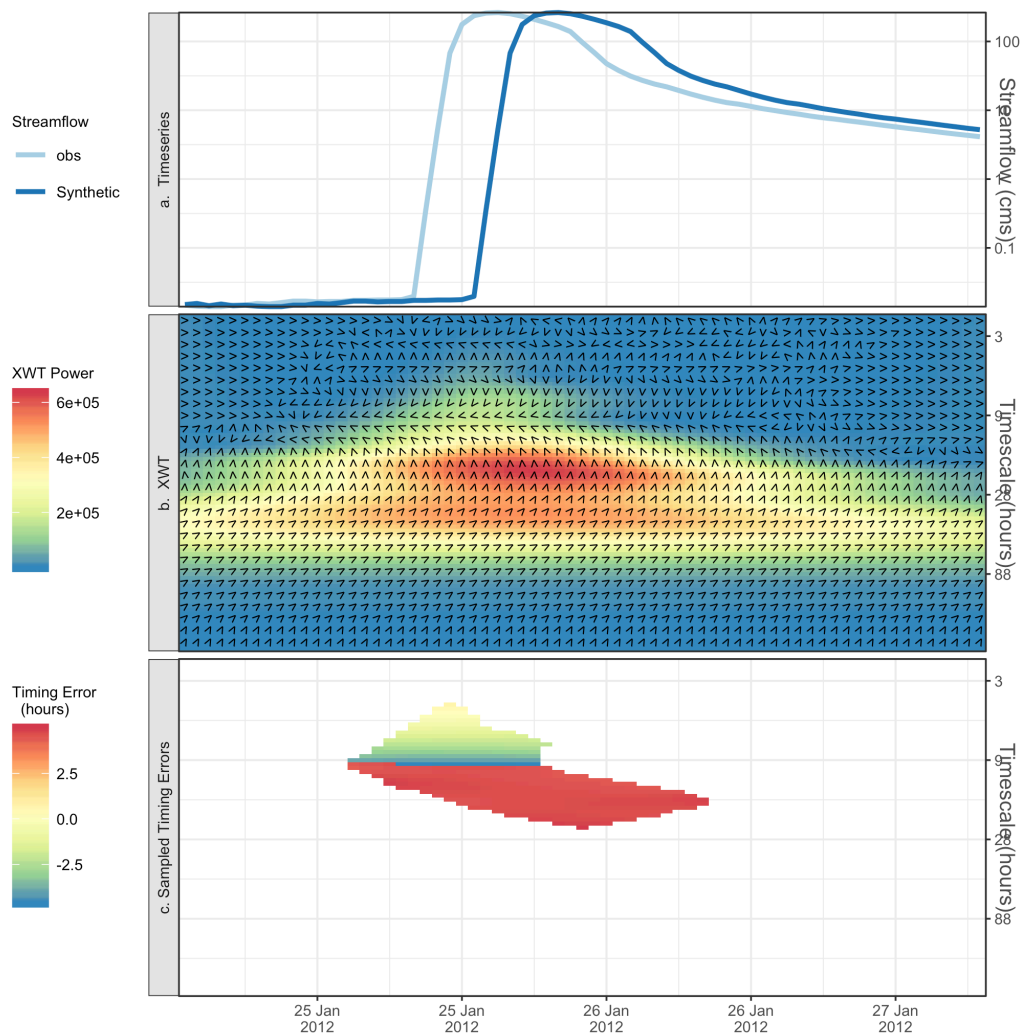


621

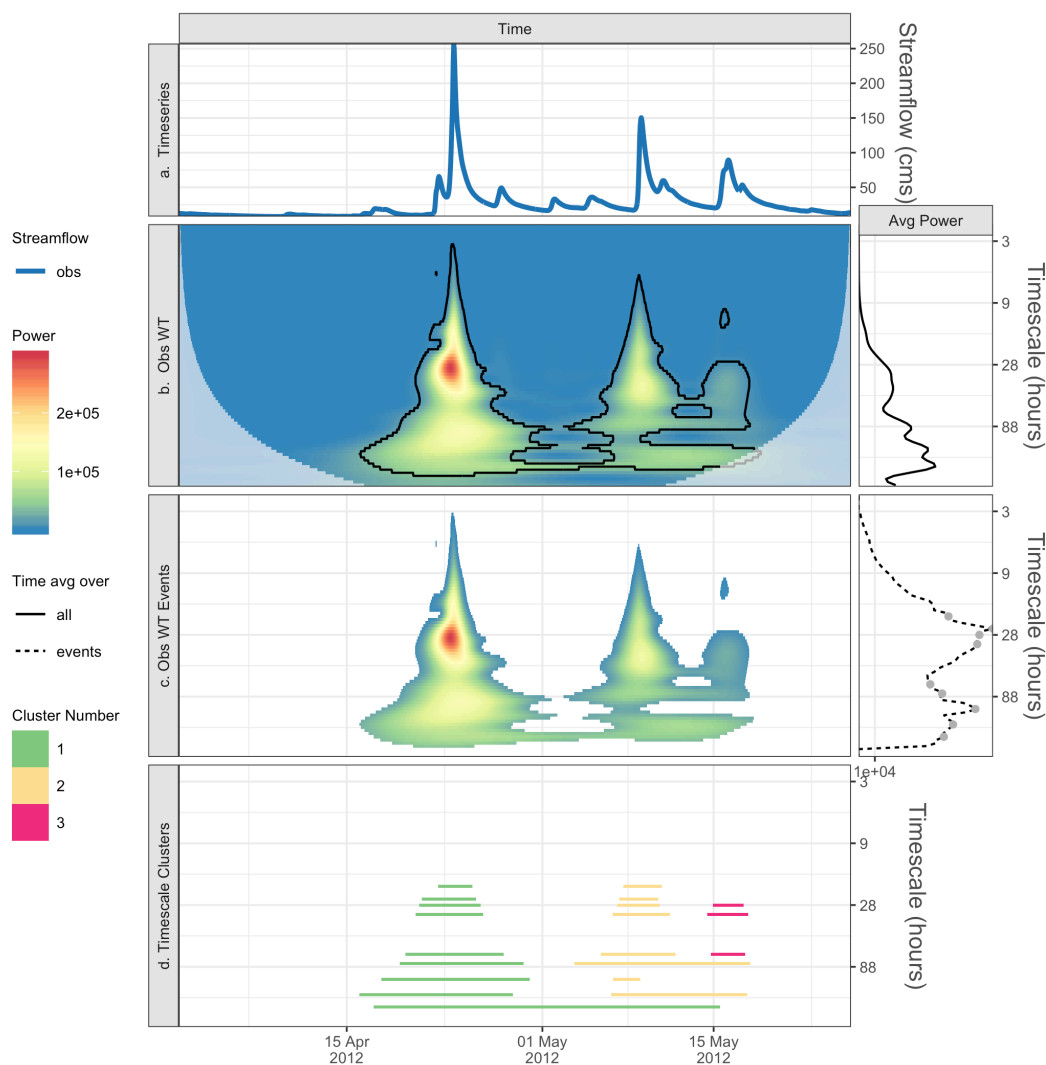
Figures



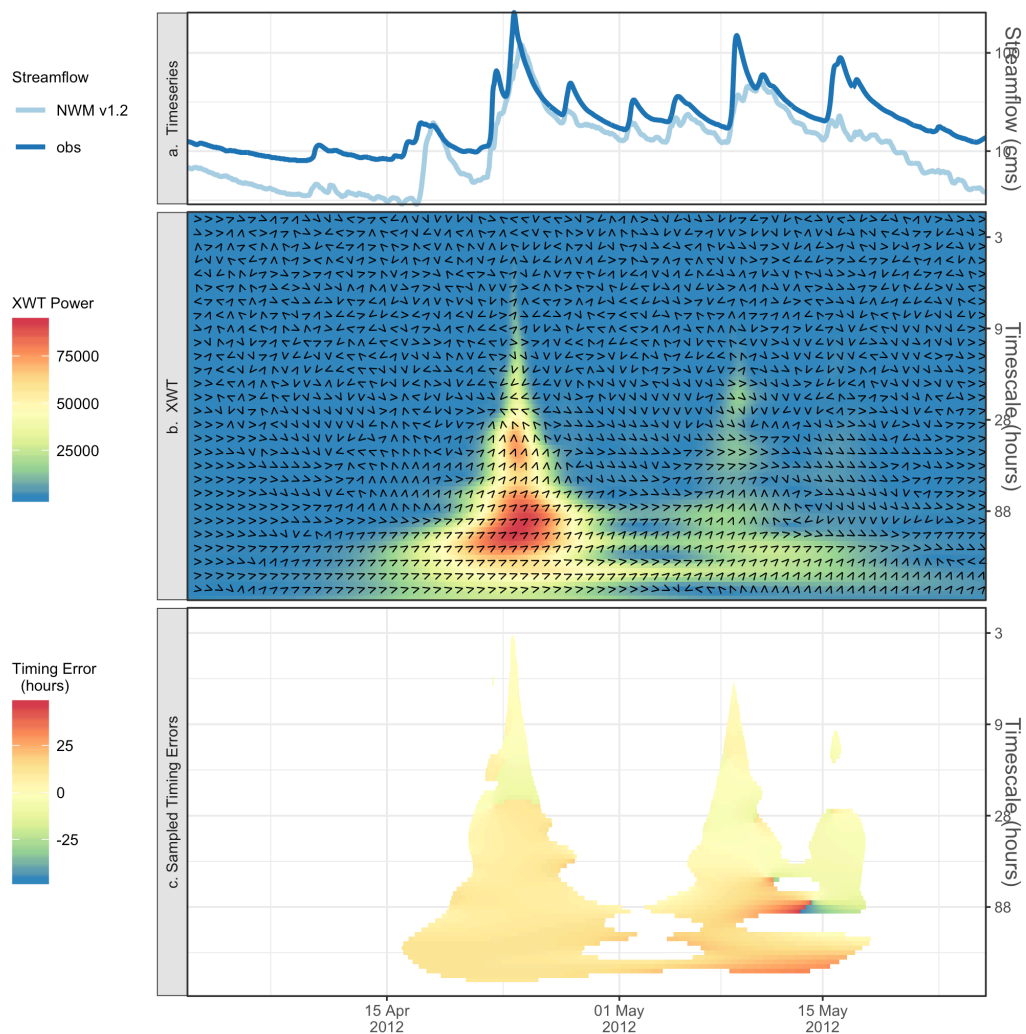
622
 623 Figure 1. An isolated peak from Onion Creek, TX: (a) observed time series, (b) observed wavelet
 624 power spectrum (left) and average power by timescale for all points (right); (c) statistically
 625 significant wavelet power spectrum or events (left) and average power by time scale for all
 626 events with maxima shown by grey dots (right); (d) Characteristic scale event cluster (horizontal
 627 green line) and maxima (star).



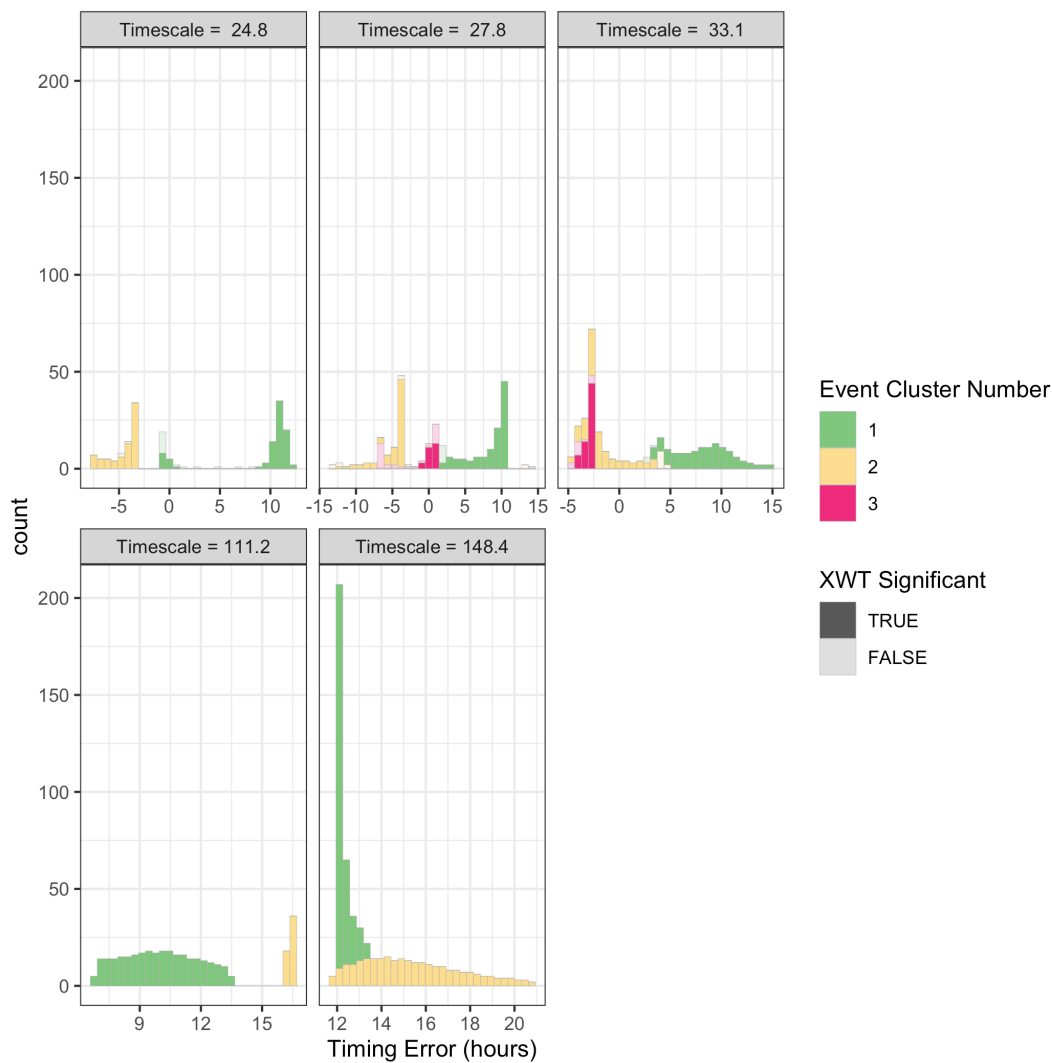
628
629 Figure 2. An isolated peak from Onion Creek, TX and a synthetic +5 hour offset: (a) observed
630 and synthetic time series, (b) cross wavelet (XWT) power spectrum and phase angles,
631 (c) sampled timing errors for observed events.



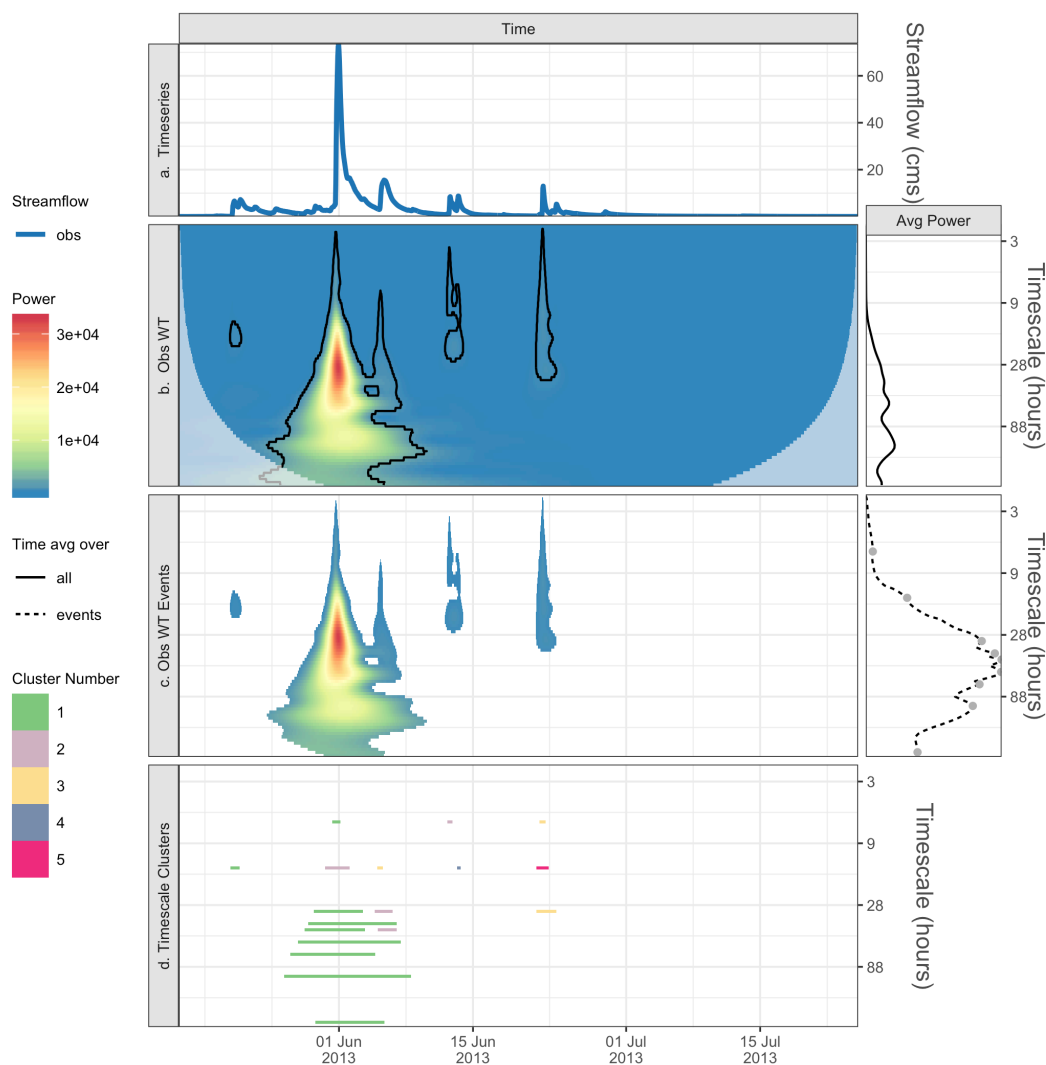
632
 633 Figure 3. Multiple peaks from Pemigewasset River, NH: (a) observed time series, (b) observed
 634 wavelet power spectrum (left) and average power by timescale for all points (right); (c)
 635 statistically significant wavelet power spectrum or events (left) and average power by time scale
 636 for all events with maxima shown by grey dots (right); (d) Characteristic scales event clusters
 637 (horizontal lines).



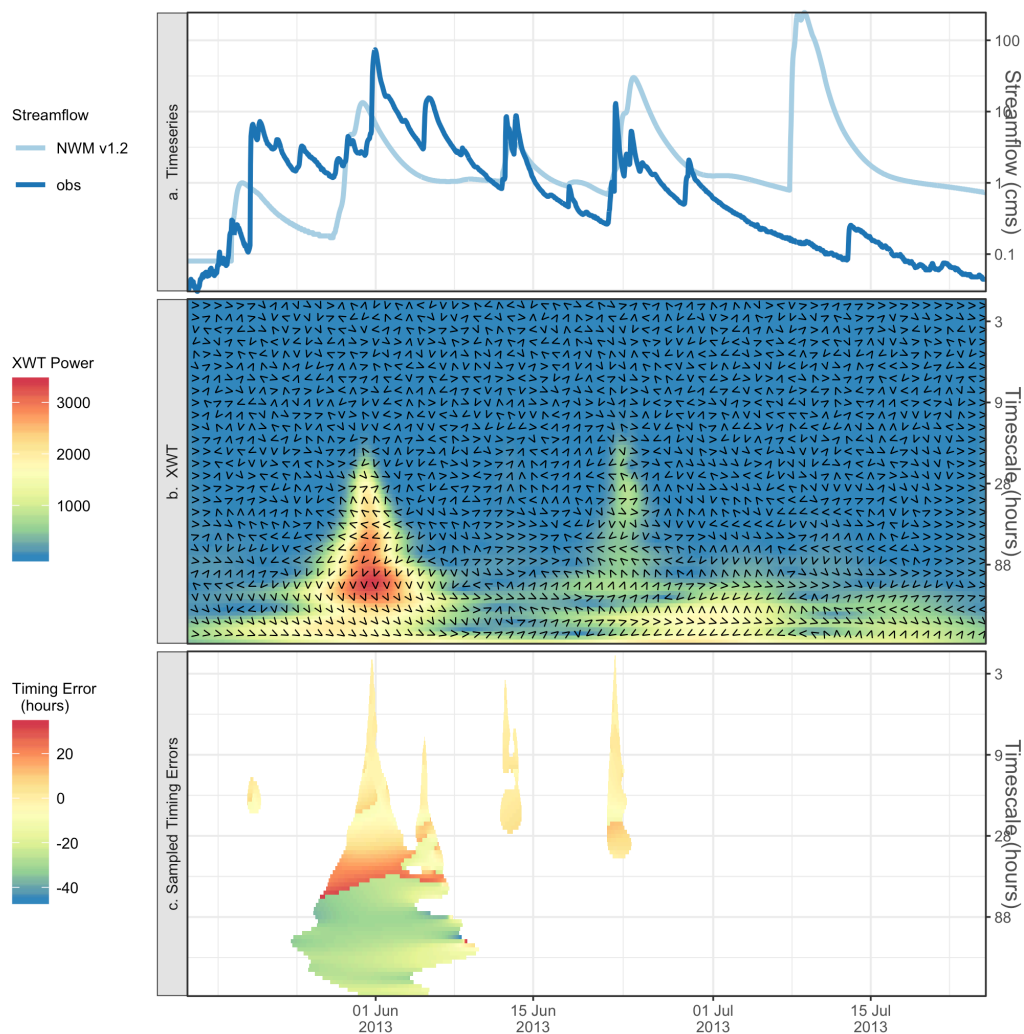
638
639 Figure 4. Multiple peaks from Pemigewasset River, NH: (a) observed and simulated NWM time
640 series, (b) cross wavelet (XWT) power spectrum and phase angles (arrows), (c) sampled timing
641 errors for observed events.



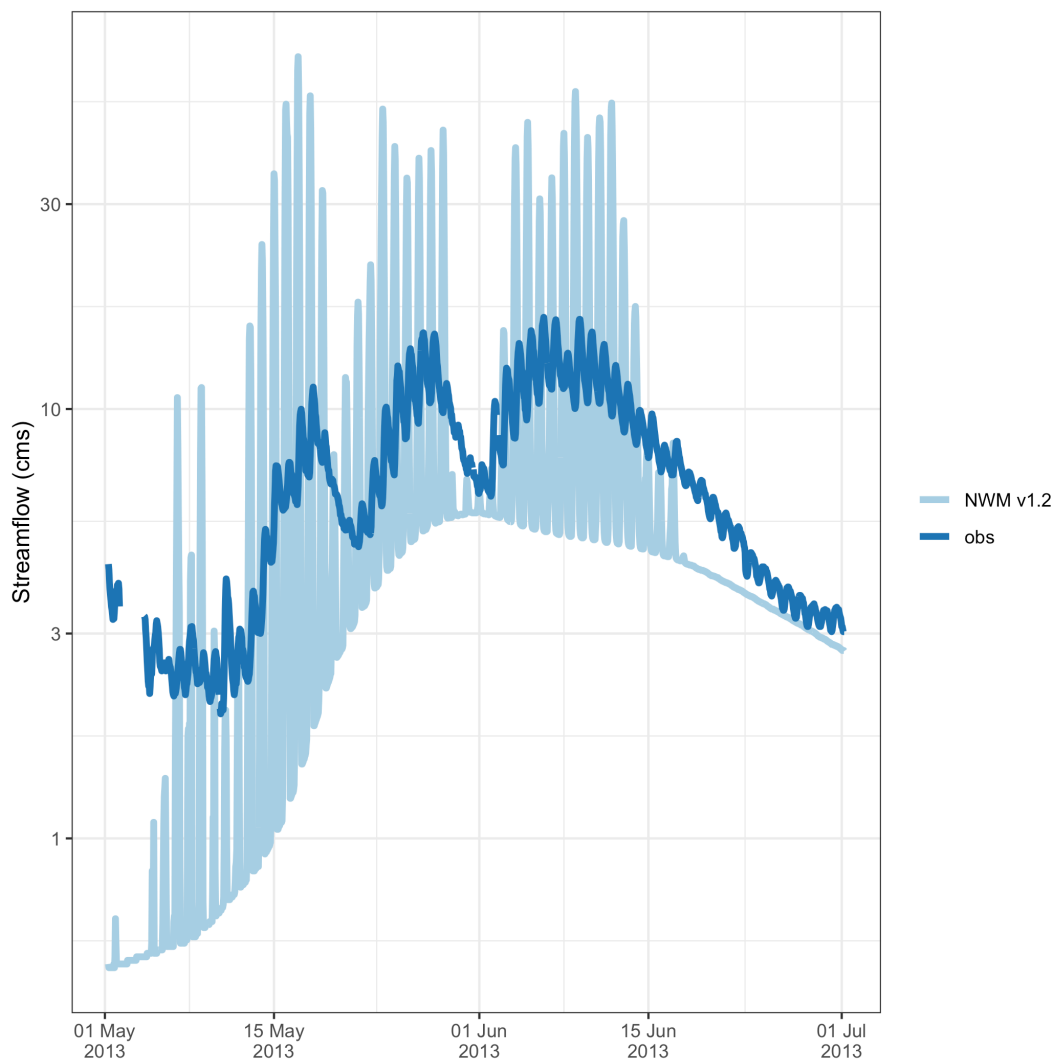
642
643 Figure 5. Multiple peaks from Pemigewasset River, NH : For the top 5 characteristic timescales
644 (see panel title), timing error distributions for event clusters. Dark colors show if the event was
645 significant in the cross wavelet transform (XWT), muted colors indicate no significance.



646
647 Figure 6. Consecutive peaks from Bad River, SD: (a) observed time series, (b) observed wavelet
648 power spectrum (left) and average power by timescale for all points (right); (c) statistically
649 significant wavelet power spectrum or events (left) and average power by time scale for all
650 events with maxima shown by grey dots (right); (d) Characteristic scales event clusters
651 (horizontal lines).

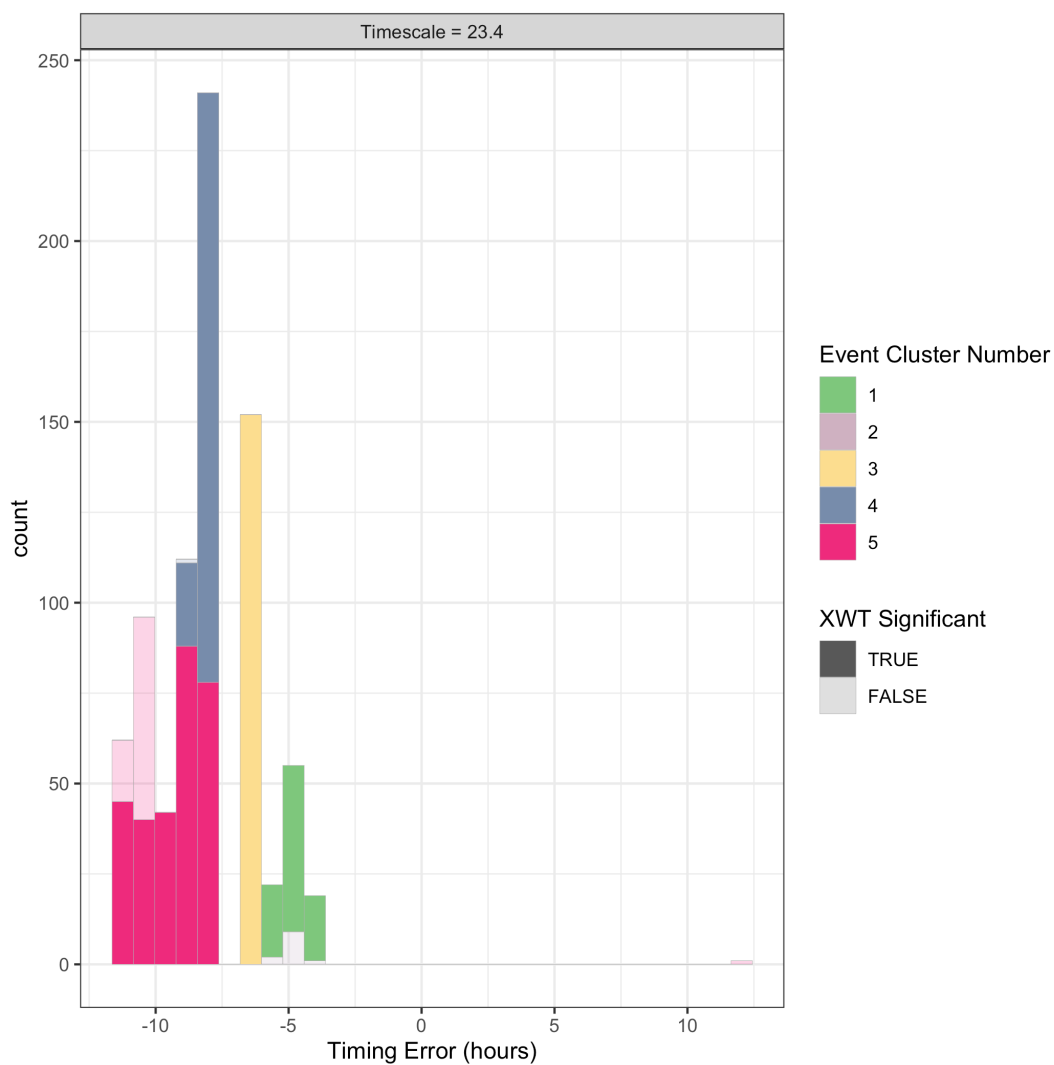


652
653 Figure 7. Consecutive peaks from Bad River, SD: (a) observed and simulated NWM time series,
654 (b) cross wavelet (XWT) power spectrum and phase angles (arrows), (c) sampled timing errors
655 for observed events.

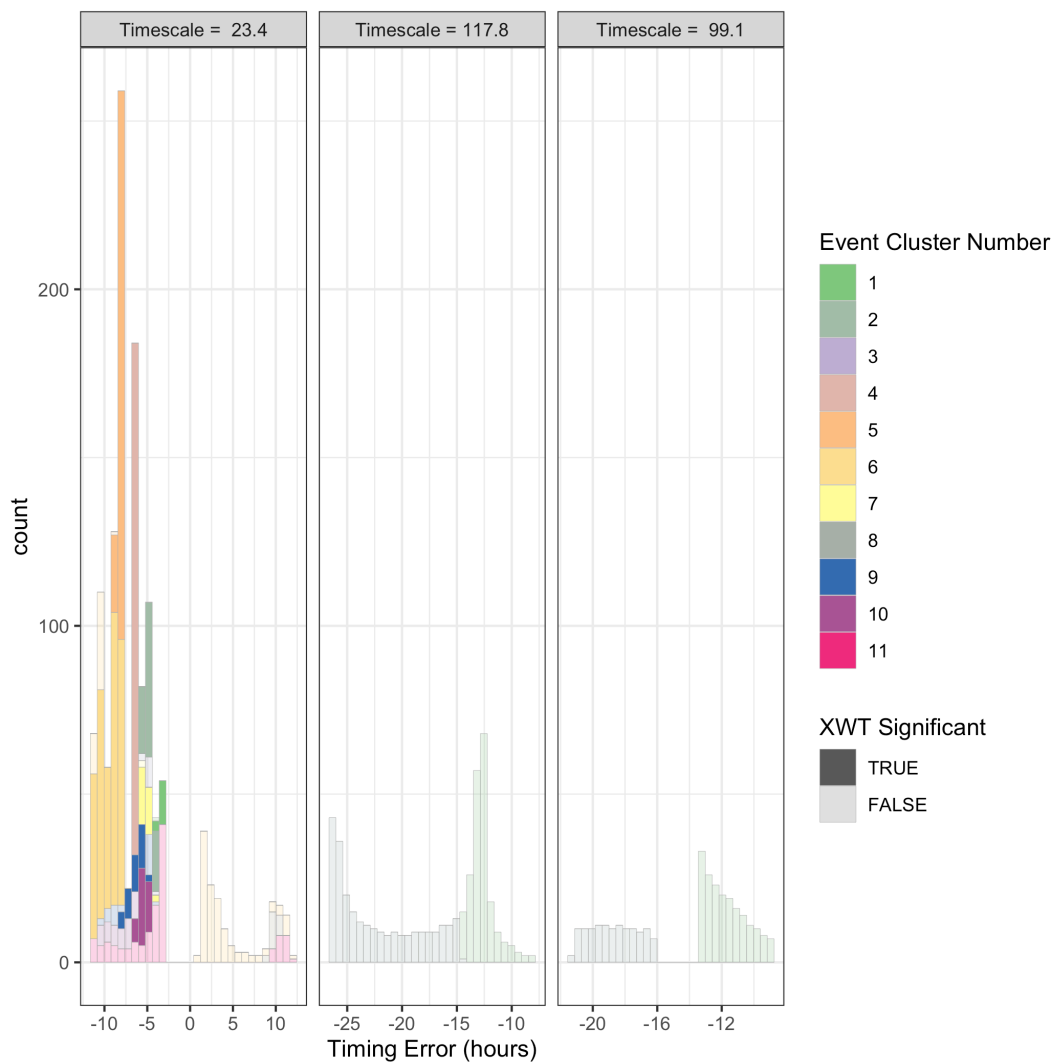


656
657

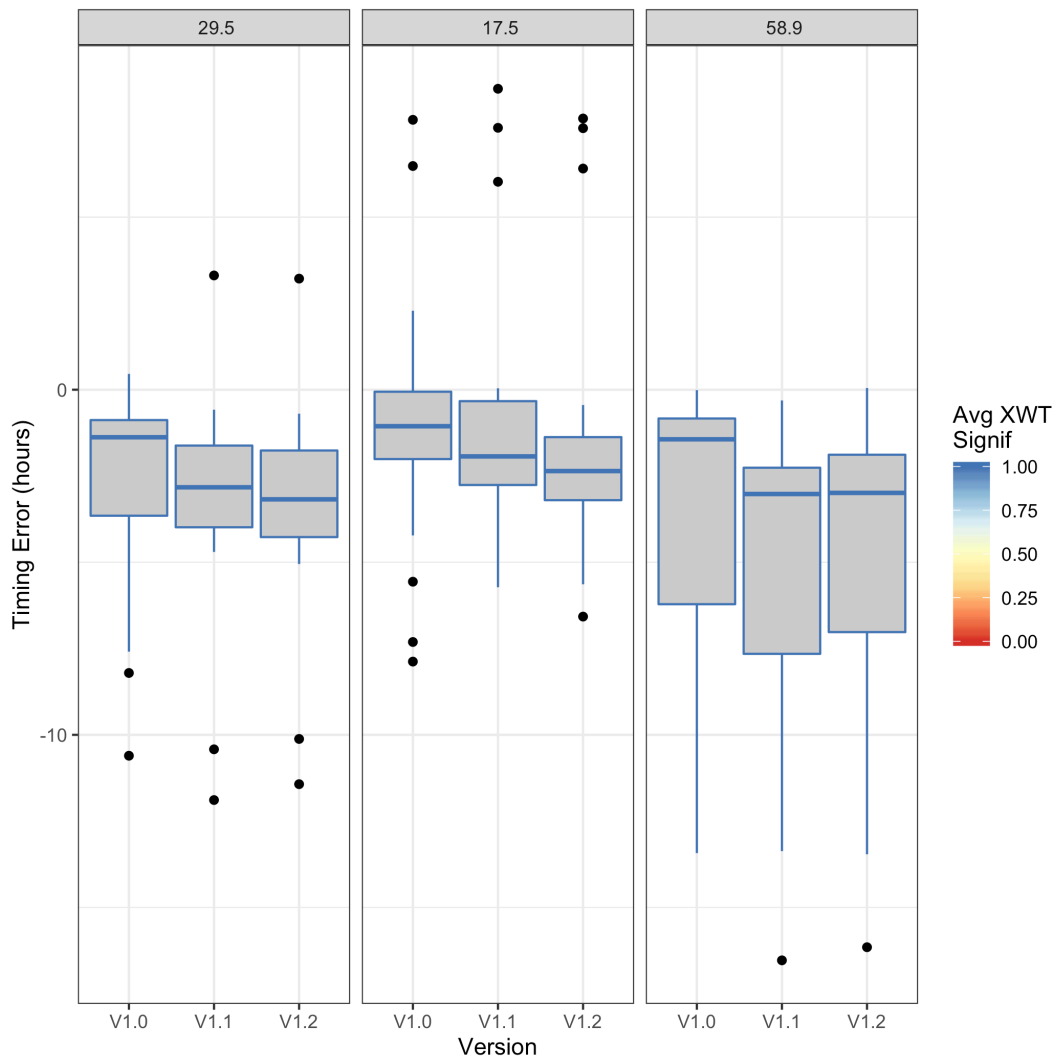
Figure 8. Taylor Park, CO: observed and simulated NWM time series.



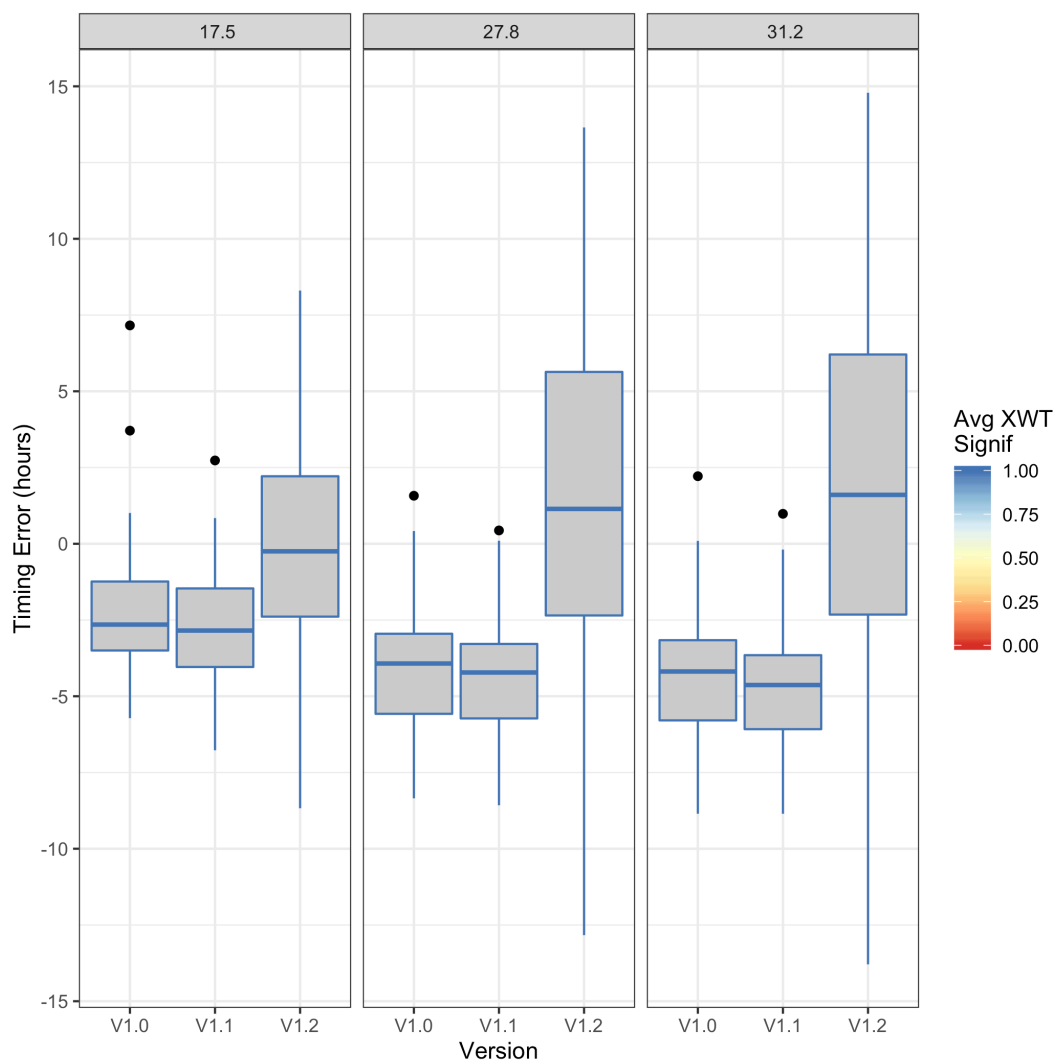
658
659 Figure 9. Taylor Park, CO: Timing error distributions for event clusters. Dark colors show if the
660 event was significant in the cross wavelet transform (XWT), muted colors indicate no
661 significance.



662
663 Figure 10. Taylor Park, CO: Timing error distributions for event clusters for top three
664 characteristic timescales (see panel title). Dark colors show if the event was significant in the
665 cross wavelet transform (XWT), muted colors indicate no significance.



666
667 Figure 11. Five year run from Onion Creek, TX: Comparing cluster max timing error
668 distributions for top three characteristic timescales (see panel title) across NWM versions;
669 outline shading shows average significance in the cross wavelet transform (XWT).



670
671 Figure 12. Five year run from Pemigewasset River, NH: Comparing cluster max timing error
672 distributions for top three characteristic timescales (see panel title) across NWM versions;
673 outline shading shows average significance in the cross wavelet transform (XWT).
674