



# A comparison of the discrete cosine and wavelet transforms for hydrologic model input data reduction

Ashley Wright<sup>1</sup>, Jeffrey P. Walker<sup>1</sup>, David E. Robertson<sup>2</sup>, and Valentijn R. N. Pauwels<sup>1</sup>

<sup>1</sup>Department of Civil Engineering, Monash University, Clayton, Victoria, Australia

<sup>2</sup>CSIRO, Land and Water, Clayton, Victoria, Australia

Correspondence to: Ashley Wright (ashley.wright@monash.edu)

Received: 19 January 2017 – Discussion started: 6 February 2017

Revised: 30 May 2017 – Accepted: 24 June 2017 – Published: 27 July 2017

**Abstract.** The treatment of input data uncertainty in hydrologic models is of crucial importance in the analysis, diagnosis and detection of model structural errors. Data reduction techniques decrease the dimensionality of input data, thus allowing modern parameter estimation algorithms to more efficiently estimate errors associated with input uncertainty and model structure. The discrete cosine transform (DCT) and discrete wavelet transform (DWT) are used to reduce the dimensionality of observed rainfall time series for the 438 catchments in the Model Parameter Estimation Experiment (MOPEX) data set. The rainfall time signals are then reconstructed and compared to the observed hyetographs using standard simulation performance summary metrics and descriptive statistics. The results convincingly demonstrate that the DWT is superior to the DCT in preserving and characterizing the observed rainfall data records. It is recommended that the DWT be used for model input data reduction in hydrology in preference over the DCT.

## 1 Introduction

Rainfall uncertainty is the biggest obstacle hydrologists face in their pursuit of accurate, precise and timely streamflow forecasts (McMillan et al., 2011). Unfortunately, errors in rainfall time series data may lead to hydrological model parameter estimates that produce adequate streamflow simulations only during the calibration period (Beven, 2006). This can lead to poor-quality streamflow predictions for independent periods and low confidence in the ability of streamflow forecasts. Consequently, a precise and accurate representation of rainfall uncertainty is paramount for robust hydrolog-

ical model parameter estimation, streamflow forecasting and quantitative precipitation forecasts (QPFs). Robertson et al. (2013) and Shrestha et al. (2015) have demonstrated that skill can be added to QPFs by postprocessing with past observations. As such, skill can be added to QPFs, and consequently flood forecasts, through developing a greater understanding of rainfall uncertainty.

The propagation of input errors in rainfall runoff modeling impedes the hydrologic community's ability to validate model structural error. Despite the vast amount of literature on rainfall measurement, estimation, statistical analysis (Testik and Gebremichael, 2010) and quality control procedures (World Meteorological Organization, 2014), a shroud of uncertainty still surrounds how rainfall and its associated uncertainty should be addressed in rainfall runoff modeling. The implementation of uncertainty analysis in many hydrological applications is also often limited by computational power.

Recent advancements in computational power as well as remote sensing have led to considerable improvements in availability and quality of hydrological observations (Cloke and Pappenberger, 2009). These improvements can be leveraged to increase the hydrological and flood forecasting knowledge base and consequently provide water policy decision makers and emergency management services with higher-quality information.

The advancement of computational power has also aided the search for hydrological model parameters that optimally simulate hydrological observations. These approaches initially focused on finding only the global optimum values of the parameters for a given objective function (Duan et al., 1994; Gan and Biftu, 1996; Thyer et al., 1999). However,

in the past two decades, it has been recognized that the uncertainties in model parameters and predictions need to be estimated. Methods that seek to estimate parameter and prediction uncertainty include Bayesian recursive parameter estimation (Thiemann et al., 2001), the limits of acceptability approach (Beven, 2006; Blazkova and Beven, 2009), the Bayesian total error analysis (BATEA) framework (Kavetski et al., 2006a, b; Kuczera et al., 2006; Thyer et al., 2009; Renard et al., 2011), the simultaneous optimization and data assimilation (SODA) (Vrugt et al., 2005), the DREAM algorithm and its variations (Vrugt et al., 2005, 2008, 2009a, b; Vrugt and Ter Braak, 2011; Laloy and Vrugt, 2012; Sadegh and Vrugt, 2014), Bayesian model averaging (Butts et al., 2004; Ajami et al., 2007; Vrugt and Robinson, 2007), the hypothetico-inductive data-based mechanistic modeling framework of Young (2013) and Bayesian data assimilation (Bulygina and Gupta, 2011). It is through the development of these parameter estimation algorithms that hydrologists are able to explore input uncertainty.

Kavetski et al. (2006b) and Vrugt et al. (2008) identified the need to represent true catchment rainfall and its associated uncertainty using parameters, both applied a parametric approach to estimating true catchment rainfall and its associated uncertainty using a rainfall multiplier to storm events. The use of a parametric representation of rainfall with an effective sampling algorithm provides the ability to jointly estimate hydrologic model parameter distributions as well as input uncertainty. As in most hydrological problems, there is a lack of sufficient data to obtain a unique solution. However, Kavetski et al. (2006b) and Vrugt et al. (2008) found there were sufficient data to estimate both hydrological model parameters and rainfall input. Data reduction transformations offer the potential to reduce the dimensionality of the parameter estimation problem and thus enable a more robust inference. Signal transforms, such as Fourier and wavelet transforms, are examples of data reduction transformations that have been applied in hydrology; however, they have not previously been used to reduce the dimensionality of input data.

Fourier transforms use sinusoidal functions to represent the spectral component of an input signal; thus, a periodic signal could be represented using a smaller number of Fourier coefficients than the number of input data points. A pitfall of the Fourier transform is that it represents the spectral components of a signal, without any indication of the time localization of those specific spectral components. In order to account for this, the windowed Fourier transform (WFT), sometimes referred to as the short-time Fourier transform, segments the signal into discrete time windows before performing the Fourier analysis. A major drawback to this approach is that the uncertainty principle of signal processing imposes a limitation on the time and frequency resolutions that can be obtained for a given signal. As a response to this, Daubechies (1990) produced discrete basis functions with good time and frequency localization. In conjunction with the pyramid algorithm, as described by Mallat (1989), this

work formed the basis for multi-resolution analysis with the discrete wavelet transform (DWT; Polikar, 1999). The DWT decomposes an input signal into high- and low-frequency components.

Wavelet analysis was first introduced to the geophysical sciences by Kumar and Foufoula-Georgiou (1997) and has been adopted for several different applications. Wavelet analysis has been used to assess the performance of hydrological models for parameter estimation (Schaeffli and Zehe, 2009) to analyze changes over different time periods for both streamflow and precipitation data (Nalley et al., 2012). Various spectral methods have also been applied in hydrology, including the application of discrete Fourier transforms to calibrate water and energy balance models (Pauwels and De Lannoy, 2011) and for the calibration of the conceptual rainfall runoff model known as the probability distributed model (PDM) (De Vleeschouwer and Pauwels, 2013). While wavelet and spectral methods have been applied in the hydrological sciences, to date there have been no instances in which the suitability of different transforms has been compared for hydrological data reduction applications. Labat (2005) has pointed out that Fourier transforms and their derivatives are not well suited to reconstruct hydrologic data, which are generated by transient mechanisms. This is due to the Fourier transforms' poor capability to represent sporadic high-frequency events when dimensionally reduced. If model input data reduction techniques are to be accepted by the hydrologic community, it is of critical importance that the transform used is able to reconstruct transient events. Through a comparative study, it will be shown that DWTs are a good multi-resolution alternative to the discrete cosine transform (DCT).

Traditionally, transform coefficients are the result of a convolution operation on an input signal. However, the aim of model input data reduction is to estimate these transform coefficients. Hence, they shall be referred to as transform parameters from herein. This paper provides novel theoretical and numerical comparisons of the DCT and DWT in a hydrological context. The ability of both transforms to reproduce key components of hydrological data sets is investigated. The extent to which each transform can reproduce hydrologic data using a decreasing number of parameters will serve as a metric upon which their ability to be used as a tool for model input data reduction for hydrological data will be evaluated. To address the requirements for hydrologic model input data reduction, this paper details (i) theoretical differences between the DCT and DWT, (ii) methodologies to reduce input rainfall to parameters and (iii) an evaluation of the proposed methodologies using several simulation performance summary metrics.

## 2 Model input data reduction theory

For this study, model input data reduction theory is introduced using a lumped conceptual watershed model. Consider a nonlinear model,  $\mathcal{F}(\cdot)$ , which simulates  $n$  discharge values,  $\widehat{\mathbf{Y}} = \{\widehat{y}_1, \dots, \widehat{y}_n\}$ , in  $\text{mm day}^{-1}$  according to

$$\widehat{\mathbf{Y}} = \mathcal{F}(\boldsymbol{\theta}, \widetilde{\mathbf{x}}_0, \widehat{\mathbf{E}}, \widehat{\mathbf{R}}), \quad (1)$$

where the model input arguments are the  $1 \times d$  vector  $\boldsymbol{\theta}$ , with arbitrary model parameter values, the  $1 \times m$  vector  $\widetilde{\mathbf{x}}_0$ , with values of the initial states in millimeters and the  $1 \times n$  vectors  $\widehat{\mathbf{E}} = \{\widehat{e}_1, \dots, \widehat{e}_n\}$  and  $\widehat{\mathbf{R}} = \{\widehat{r}_1, \dots, \widehat{r}_n\}$  which store the observed values of the potential evapotranspiration (PET) and rainfall in  $\text{mm day}^{-1}$ , respectively. Note that  $\widehat{\mathbf{R}}$  is used to represent rainfall and not precipitation, as snow, hail and other forms of precipitation are not considered. The  $\widehat{\phantom{x}}$  (hat) symbol is used to denote measured quantities and the  $\widetilde{\phantom{x}}$  (tilde) symbol reflects variables that are either reconstructed or could, in theory, be observed in the field but due to their conceptual nature are difficult to determine accurately.

If the traditional hydrological perspective in which the inputs  $\mathbf{E}$  and  $\mathbf{R}$  are considered to be fixed and known quantities is relaxed, and rainfall is now considered unknown, then a new inference problem arises in which the input rainfall is estimated via the treatment of the input rainfall as a series of parameters. Inference problems in which the input is considered unknown can be dealt with using a Bayesian framework. Such inference problems have been considered by Kavetski et al. (2006a) and Vrugt et al. (2008) but are outside the scope of this paper. Consequently, for rainfall to be inferred, a suitable parametric representation of rainfall must be determined.

Given a daily rainfall data record with  $n$  observations in millimeters,  $n$  rainfall parameters could be used to represent the input hyetograph. This approach would be particularly elegant and parsimonious. Yet, for a 10-year record of daily discharge data, the inference problem would grow from  $d$  model parameters to roughly  $10 \times 365 + d = 3650 + d$  parameters. These values would need to be estimated from the observed rainfall and discharge data record, respectively. As many hydrological models are already underdetermined, the introduction of additional parameters would make the model even less determinable. Additionally, an excessive amount of CPU time is required to solve for a 3600+ dimensional posterior parameter distribution. An alternative approach is therefore necessary.

Sparse transforms convey large amounts of data using fewer parameters than data points in the observed signal. An input rainfall signal can be reduced to sparse transform parameters. Doing so allows multiple rainfall observations to be modified using a single parameter. Some or all of these transform parameters can be altered before the transform is inverted to produce a new input signal for streamflow simulation and posterior analysis. The use of sparse transforms to represent input time series enables input uncertainty to be

explored in great detail. The ability of discrete wavelet and Fourier transformations to reduce hydrological input data to a set of parameters for uncertainty estimation is compared using theoretical and analytical methods.

### 2.1 Overview of the DCT and DWT

Wavelet and Fourier transforms are invertible transforms in which a forward convolution operation can be used to decompose a signal into various components. Similarly, a backwards deconvolution operation can be applied to retrieve the original signal. Fourier-based transforms decompose signals into frequency components and are best used for regular time-invariant signals that do not exhibit time-specific information. Alternatively, wavelet-based transforms decompose signals into frequency and time components. The advantage of using wavelet functions to transform data is that time-specific information about when higher frequency components occur can be preserved. To obtain time-specific information, Fourier-based transforms can be applied over pre-specified temporal windows. Yet, this approach is limited by the uncertainty principle of signal processing. The uncertainty principle of signal processing imposes a lower limit on obtainable resolutions in the time–frequency domain such that

$$\sigma_t \sigma_\omega \geq \frac{1}{2}, \quad (2)$$

where  $\sigma_t$  (s) and  $\sigma_\omega$  ( $\text{s}^{-1}$ ) are the respective temporal and frequency widths used in the sparse transform.

Applying the uncertainty principle of signal processing (Eq. 2), it is clear that any attempt to narrow the temporal period analyzed to gain increased resolution in the time domain would be met by a widening of the frequency spectrum and consequently a loss of resolution in the frequency domain.

Considering that there is no time–frequency window that is able to obtain limitless resolution in both the time and frequency domains, it is clear that an alternative solution must be found. Wavelet transforms can be used to decompose a signal into different levels that consist of different time and frequency resolution windows. Thus, the wavelet transform is able to be configured to simultaneously obtain high levels of resolution in both the time and frequency domains. For a more detailed discussion on wavelets and sparse transforms, the reader is referred to Mallat (2009).

### 2.2 Discrete cosine transform

The DCT (Ahmed et al., 1974) is a version of the WFT that has advantageous properties for the field of data compression. Due to the boundary conditions of the cosine function, the DCT is well suited to represent an observed input signal with a minimal number of parameters: in this case, rainfall

$\widehat{\mathbf{R}}(t)$ . The DCT parameters  $\mathbf{p}(i)$  are calculated as

$$\mathbf{p}(i) = w(i) \sum_{t=1}^n \widehat{\mathbf{R}}(t) \cos \left[ \frac{\pi}{2n} (2t - 1)(i - 1) \right], \quad (3)$$

where  $i = 1, 2, \dots, n$  and

$$w(i) = \begin{cases} \frac{1}{\sqrt{n}}, & i = 1 \\ \sqrt{\frac{2}{n}}, & 2 \leq i \leq n. \end{cases} \quad (4)$$

The convolution process can be reversed to reconstruct the observed signal using the inverse transform:

$$\widetilde{\mathbf{R}}(t) = \sum_{i=1}^n w(i) \mathbf{p}(i) \cos \left[ \frac{\pi(2t - 1)(i - 1)}{2n} \right], \quad (5)$$

where  $t = 1, 2, \dots, n$ .

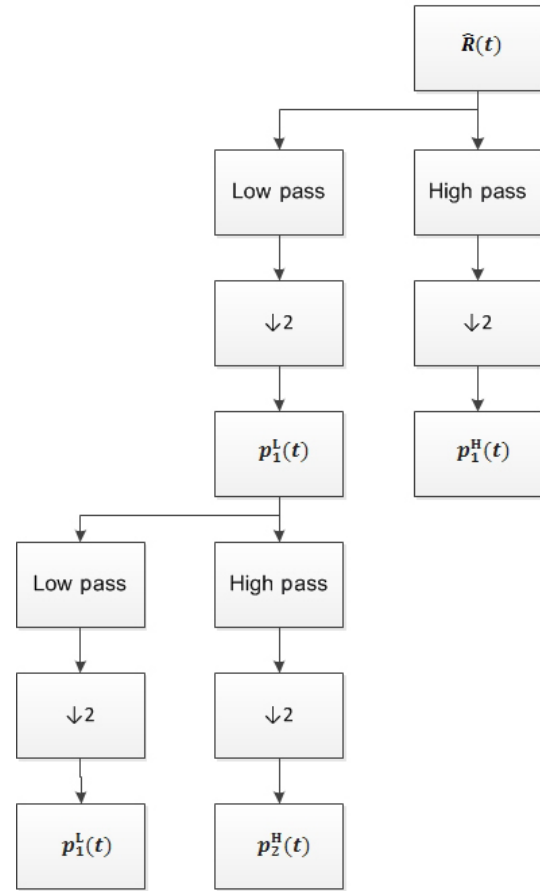
### 2.3 Discrete wavelet transform

Using the pyramid algorithm, depicted in Fig. 1, Mallat (1989) first described the decomposition of an input signal into multi-resolution components using high- and low-pass filters. Each stage of decomposition is referred to as a level. An advantage of using wavelets is that decomposition can be performed using a variety of different wavelet families. This allows for signals with differing properties to be analyzed using the same methodology. The most commonly used wavelet family is the Daubechies wavelets (Daubechies, 1990). Each wavelet within each family consists of a scaling  $h(m)$  and wavelet  $w(m)$  function, where  $m$  denotes the length along the scaling and wavelet function. The scaling and wavelet functions are used in the low- and high-pass filtering sequences, respectively. Whilst there are numerous wavelet families that can be chosen for analysis, this study applies the most commonly used Daubechies wavelets. Depending on the choice of wavelet, stepwise convolutions of the input signal are performed over the filter length  $L$ .  $j_{\max}$  imposes an upper limit on the level of decomposition  $j$  that a signal can be decomposed into, where

$$j_{\max} = \left\lfloor \log_2 \left( \frac{n + L - 1}{2} \right) \right\rfloor, \quad (6)$$

in which  $\lfloor \cdot \rfloor$  is the floor operator. The input signal is then convoluted by being passed through high- and low-pass filters, where

$$\mathbf{p}_j^L(i) = \begin{cases} \sum_{m=1}^L \widetilde{\mathbf{R}}(2i - m - 1)w(m), & j = 1 \\ \sum_{m=1}^L \mathbf{p}_{j-1}^L(2i - m - 1)w(m), & j > 1 \end{cases} \quad (7)$$



**Figure 1.** A schematic showing the pyramid algorithm used to decompose and downsample ( $\downarrow 2$ ) an input signal ( $\widehat{\mathbf{R}}$ ) into high- and low-frequency components. The input signal is filtered using the high- and low-pass filters described in Eqs. (7) and (8) before being downsampled to produce the level 1 high- and low-pass parameters. The low pass parameters are now used as input for the high- and low-pass filters. This process of filtering and downsampling is repeated until the desired level of decomposition is met.

is the low pass and

$$\mathbf{p}_j^H(i) = \begin{cases} \sum_{m=1}^L \widetilde{\mathbf{R}}(2i - m - 1)h(m), & j = 1 \\ \sum_{m=1}^L \mathbf{p}_{j-1}^L(2i - m - 1)h(m), & j > 1 \end{cases} \quad (8)$$

is the high pass,  $i = 1, \dots, n_{j-1} + L - 1$  and refers to the  $i$ th parameter,  $j = 1, \dots, j_{\max}$  and refers to the  $j$ th level,  $m$  refers to the  $m$ th filter coefficient. The resultant low-pass  $\mathbf{p}_j^L(i)$  and high-pass  $\mathbf{p}_j^H(i)$  parameters are commonly referred to as approximation and detail parameters, respectively. After the input signal is passed through the high- and low-pass filters there is an issue of redundancy that needs to be dealt with. The filters split the input signal into high- and low-frequency components that each contain roughly half the information

of the input signal. As the length of each of the resultant approximation and detail parameter series is equivalent to the length of the input signal, each of the parameter series must be downsampled. The process of downsampling removes every other parameter. It is the process of high- and low-pass filtering followed by downsampling that enables the DWT to analyze multi-resolution components of a signal. After downsampling, the length of the resultant approximation and detail parameter series is

$$n_j = \begin{cases} \left\lfloor \frac{n+L-1}{2} \right\rfloor, & j = 1 \\ \left\lfloor \frac{n_{j-1}+L-1}{2} \right\rfloor, & j > 1 \end{cases}, \quad (9)$$

where  $n_j$  refers to the length of the series at the  $j$ th level. If further decomposition is required, the downsampled low pass may be fed back into the filters until the resultant parameters can no longer be split any further. An iteration of this process is shown in Fig. 1. To reverse the decomposition process and reconstruct a signal, upsampling is performed on the parameter series before the lower level parameters are obtained through

$$p_{j-1}(i) = \sum_{m=\lceil i/2 \rceil}^{\lceil (L-1+i)/2 \rceil} \left( p_j^H(i)h(2m-i) \right) \left( p_j(i)w(2m-i) \right), \quad j > 1, \quad (10)$$

where  $\lceil \cdot \rceil$  is the ceiling operator and the input signal is reconstructed using

$$\tilde{R}(i) = \sum_{m=\lceil i/2 \rceil}^{\lceil (L-1+i)/2 \rceil} \left( p_j^H(i)h(2m-i) \right) \left( p_j(i)w(2m-i) \right), \quad j = 1. \quad (11)$$

### 3 Data

This study utilizes data from the Model Parameter Estimation Experiment (MOPEX) data set. The 10 years of rainfall data spanning the 1990s for 438 catchments in the United States of America (USA) are used to compare the suitability of the DWT and DCT to represent rainfall time series. The catchments used in this study were chosen to ensure they had sufficient rain gauge density and represented a range of catchment sizes and climates. Rainfall for the Leaf River catchment (Collins, Mississippi), a catchment that is frequently used for hydrological studies (Sivakumar, 2001; Tang et al., 2006; Bulygina and Gupta, 2011), is used to compare the DWTs' and DCTs' ability to reconstruct high-magnitude rainfall events. A single rainfall product for each catchment is used for analysis at a daily time step. A complete description of the selection process and MOPEX data set is given by Schaake et al. (2006). No streamflow data are used in the experiment.

### 4 Experiment design

This experiment does not involve the use of any hydrological models. Due to this and the nature of the transforms, there are no calibration and evaluation periods. A major use of both the DWT and DCTs has been in image compression; consequently, the observed input signals were compressed and decompressed using a methodology similar to that used in image compression. In order to determine which transform's parameters are able to effectively store the most hydrological input data, both DWT and DCT parameters will be compressed to varying extents for the MOPEX rainfall time series.

The process undertaken involves a number of steps. Firstly, before any compression is applied, the original rainfall signal for a given catchment is transformed into DCT and DWT parameters using Eqs. (3) and (4) and Eqs. (7) to (9) for the DCT and DWT, respectively. Secondly, each transform is compressed by iteratively zeroing out parameters that provide a low degree of information; these parameters are those closest to zero. A threshold value  $T$  (mm) applies a lower limit for which transform parameters above the threshold are retained. This threshold is iteratively increased until the compressed transform is composed of the desired number of remaining parameters  $k$  and percent of original parameters (POP) is met.

$$\text{POP}(T) = 100 \times \left( \frac{k}{n} \right), \quad (12)$$

where  $k$  becomes smaller as the threshold  $T$  increases and  $\lim_{T \rightarrow \infty} \text{POP} = 0$ . The next step is to reconstruct the observed signal from the compressed transform parameters using Eqs. (5) and (11) for the DCT and DWT, respectively. After the reconstruction has been performed, a comparison between the reconstructed and observed rainfall can be made. Lastly, this process is iterated for different POPs as well as for each catchment within the data set.

To provide a meaningful comparison between the DCTs' and DWTs' ability to reproduce different rainfall time series with an increasing POP, a number of simulation performance summary metrics are used. Following Moriasi et al. (2007), a combination of graphical techniques and dimensionless and error index statistics that are widely accepted by the hydrological community were adopted for model evaluation. The Nash–Sutcliffe efficiency (NSE) and the root mean square error (RMSE) to standard deviation ratio (RSR) of the observed input signal ( $\text{RSR} = \text{RMSE}/\sigma_{\text{obs}}$ ) are used to compare the performance of the reconstructed rainfall signal with the observed rainfall signal. Once the reconstructed signals are obtained, further comparison with the observed rainfall will be made using the bias summary metric. The variance, kurtosis and skewness of the reconstructed signals will be compared with those of the observed signal. The bias is calculated as  $\sum_{t=1}^n [\hat{R}(t) - \tilde{R}(t)]/n$ , where  $\hat{R}(t)$  and  $\tilde{R}(t)$  are the observed and reconstructed rainfall signals, respectively. The recon-

structured variance, kurtosis and skewness are all normalized by the observed input signals variance, kurtosis and skewness, respectively. The peak error (PE) is the peak rainfall error over the 10-year period. It is used to compare the reconstructed and observed signals for seasonal and flood forecasting situations. The PE is normalized by the peak height of the observed input signal. Further, the number of rain events missed is computed for each reconstruction by flagging original or reconstructed observations that exhibit no rainfall. Either the absolute difference between the reconstructed and original observation is less than 0.01 or the ratio of the reconstructed and original observation is equal to 0 or larger than 10. Lastly, reconstructed rainfall using the DCT and DWT will be presented for the Leaf River catchment to compare each transform's ability to reconstruct high-magnitude rainfall events.

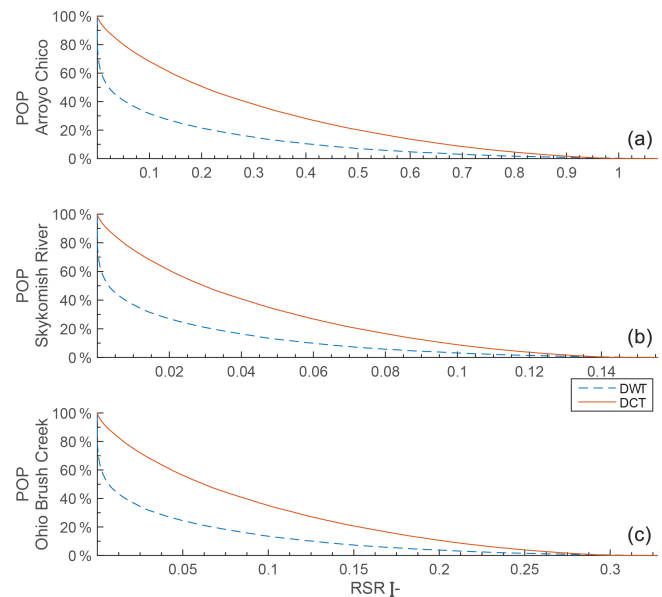
## 5 Results

Figure 2 shows the relationships between RSR and the number of transform parameters using the DCT and DWT for three different catchments: Arroyo Chico, Skykomish River and Ohio Brush Creek. These catchments represent the smallest, largest and mean rainfall volumes for the MOPEX data set, respectively. It is clear that for all but the highest POP the DWT is able to reconstruct the observed signal with lower RSR than the DCT and that as the rainfall volume increases the RSR decreases. For intermediate POPs, the DWT is able to reconstruct the observed signal with significantly better RSR than the DCT. As the POP approaches both 100 and 0%, there is little discernible difference between the DCT and DWT reconstructions.

By comparing the reconstructed DWT and DCT signals, using 20 POP and the observed rainfall signal as a reference, a histogram for the NSE is shown for all catchments in Fig. 3. Each frequency count in the histogram represents a catchment from the MOPEX data set. The reconstructed DWT signals are clearly able to better simulate the observed rainfall signal. All DWT reconstructed rainfall signals obtained a higher NSE than the DCT reconstructed rainfall signals. Table 1 shows that as the transforms are compressed and fewer parameters are used in the reconstruction, the mean NSE for the DWT stays much closer to the ideal value of 1 than the DCT. Further, the standard deviation of NSE becomes much larger for the DWT.

Figure 4 compares the RSR for the DCT and DWT using four different POPs. A 1 : 1 line is included in all subplots and each point represents a catchment from the data set. If the data points fall above the 1 : 1 line, then for that catchment and POP the DWT is able to reconstruct the input rainfall signal with lower RSR. Again, it is found that DWT is always able to reconstruct the original signal with lower RSR than the DCT reconstructions for all POPs. In a similar fashion to that discussed regarding Fig. 2, it is observed that as the POP

## A. Wright et al.: Model input data reduction in hydrology



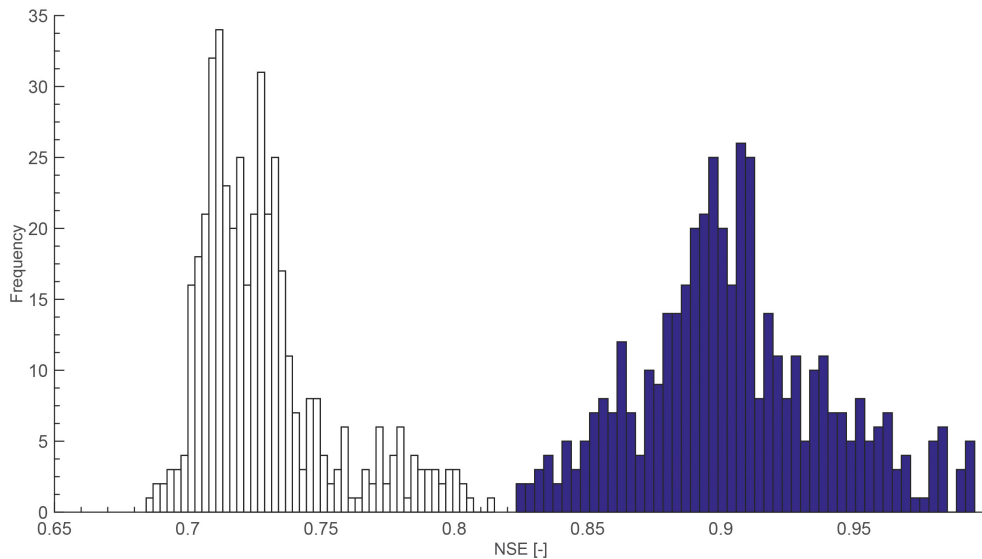
**Figure 2.** Empirical plots showing the relationship between RSR and the POP used for reconstructing an input rainfall signal using the DWT and DCT. The three catchments, from the top to the bottom of the figure, represent the smallest, largest and mean rainfall volumes throughout the 1990s for the MOPEX data set.

**Table 1.** The mean and standard deviation (SD) of NSE for the DWT and DCT using a different POP.

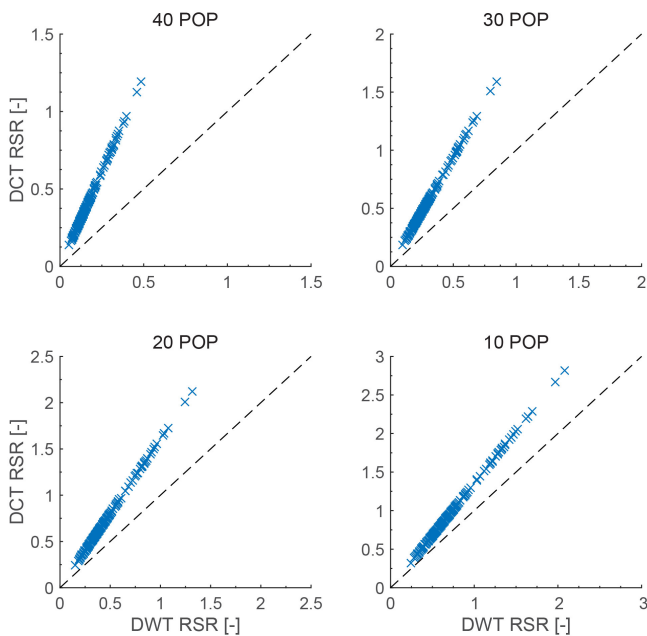
POP	NSE DWT		NSE DCT	
	Mean	SD	Mean	SD
40 %	0.988	0.007	0.918	0.010
30 %	0.965	0.017	0.844	0.016
20 %	0.905	0.036	0.729	0.025
10 %	0.746	0.070	0.522	0.037

approaches 0% the difference between the DWT and DCT reconstructions becomes smaller.

The bias, variance and skewness observed in the reconstructed signals for each catchment are shown in Fig. 5 for different POPs. The DWT reconstructions are able to maintain a smaller bias than the DCT reconstructions at different POPs for all of the catchments. As the POP decreases, the bias becomes increasingly positive and negative for the DWT and DCT, respectively. The distribution of the bias becomes more dispersed for both the DCT and DWT as the POP decreases. The bias can be seen to be dependent on the transform and POP used as well as the catchment being analyzed. Both the DWT and DCT never reconstruct the observed signal with greater variance than that of the observed rainfall signal. As the POP decreases, the normalized variance for the DCT moves further away from unity than the normalized variance for the DWT. The reduction in normalized variance means that, as the POP decreases, both the DWT and espe-



**Figure 3.** Histogram representing the reconstructed DWT (dark bins) and DCT (clear bins) NSE when compared to the observed rainfall signal. Rainfall is reconstructed after the input signal is compressed to 20 POP. Each frequency count represents a catchment from the MOPEX data set.

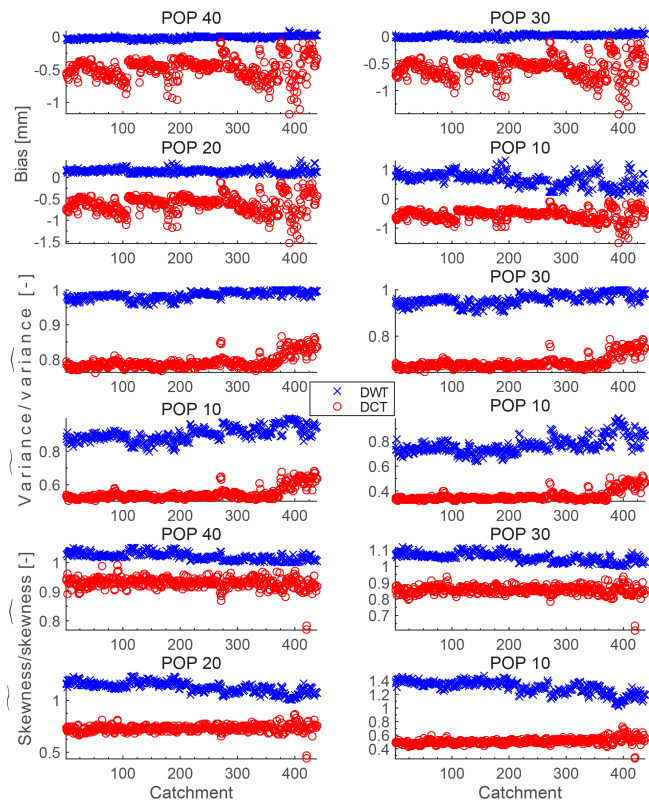


**Figure 4.** Comparative plots of RSR for the DCT and DWT using a different POP. Each data point represents a catchment.

cially the DCT reconstructions will have fewer extreme values when compared to the observed rainfall. The normalized skewness is a measure of symmetry that describes whether or not the reconstructed signal is more positively skewed (more than 1) or less positively skewed (less than 1) than the observed input signal. All of the reconstructed and observed signals had a positive normalized skewness. When

compared to the observed signal, the DWT becomes increasingly skewed as the POP is reduced. The opposite of this is observed for the DCT. This indicates that, when compressed, the DWT and DCT will reconstruct the observed rainfall signal with a greater and lower number of values close to zero when compared to the observed signal, respectively. This does not mean that the total volume will be any lower than the total volume of rainfall observed. This is made evident by the low bias observed in Fig. 5.

The normalized kurtosis and PE for all catchments using different POPs are shown in Fig. 6. The measure of kurtosis describes how much the fraction of the distributions' variance is explained by extreme deviations. Consequently, a normalized kurtosis value larger than 1 indicates that the reconstructed signals variance is explained more by extreme deviations than the observed input signal. This is likely to be the result of more rainfall values being reconstructed at the extremities than those of the observed rainfall series. A value smaller than one indicates that the variance is described less by extreme deviations than the observed input signal. Similarly, this is likely to be the result of fewer rainfall values being reconstructed at the extremities than those of the observed rainfall series. It is worth noting that a reconstructed time series can have the same variance yet different kurtosis than the observed rainfall time series. As the POP decreases, the dispersion of normalized kurtosis and skewness increases, and the normalized kurtosis and skewness for the DWT and DCT reconstructions become larger and smaller than unity, respectively. With decreasing POP, the normalized PE for the reconstructed DWT signal remains small and relatively consistent when compared to the normalized PE for the reconstructed DCT signal.

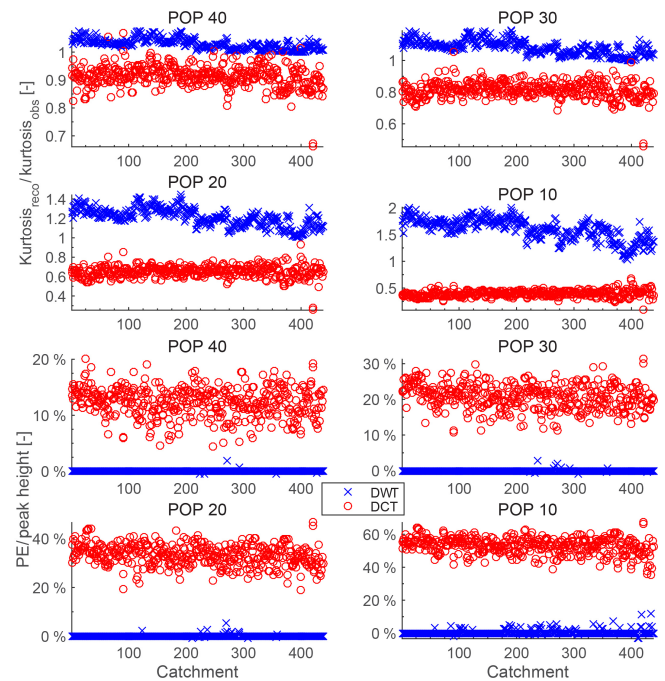


**Figure 5.** Bias and normalized variance and skewness of the reconstructed DWT and DCT signals for each catchment using a different POP.

## 6 Discussion

Figure 3 shows that the DWT and DCT are able to reconstruct the observed input signals with good efficiency using 20 POP. However, the DWT consistently outperforms the DCT. Fig. 2 shows that, as the POP decreases from 100 %, the DWT is able to reconstruct the input signal with increasingly lower RSR than the DCT; the gap in performance is largest for 40 POP. As the POP continues to decrease towards 0 %, the gap in RSR reduces to zero. It is interesting to note that the DWT perfectly reconstructs the observed rainfall signal with as many parameters as there are rainy days, whereas the DCT does not.

As the bias for the DWT is consistently close to zero, the use of the DWT for rainfall input data reduction is likely to be beneficial for hydrologic studies that have short time steps and involve rainfall as an input. Whilst modification of the DWT parameters may slightly overestimate input rainfall, it is not as significant as the consistent underestimation of input rainfall by the DCT. The diminishing ability of both the DWT and DCT to match the input rainfall signal variance indicates that both transforms smooth out input data towards the mean. This behavior is more significant for the DCT than the DWT. Consequently, when used as a technique for input data reduc-

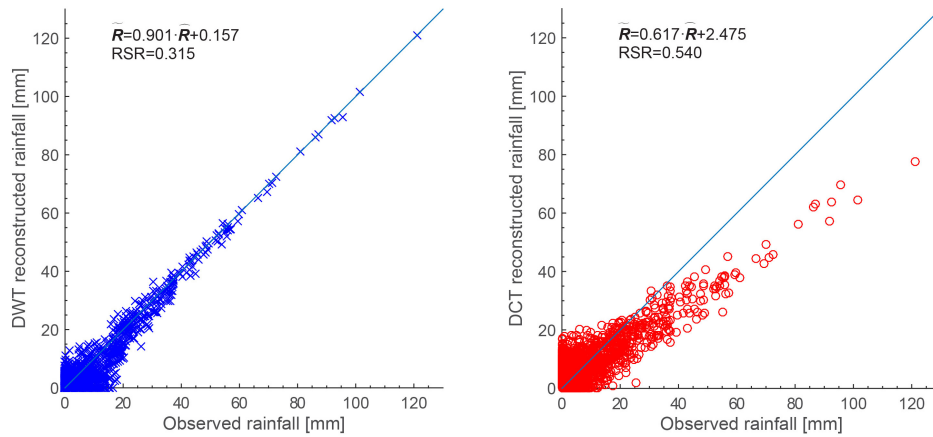


**Figure 6.** Normalized kurtosis of the reconstructed DWT and DCT signals and percentage PE for the reconstructed DWT and DCT signals for each catchment using a different POP.

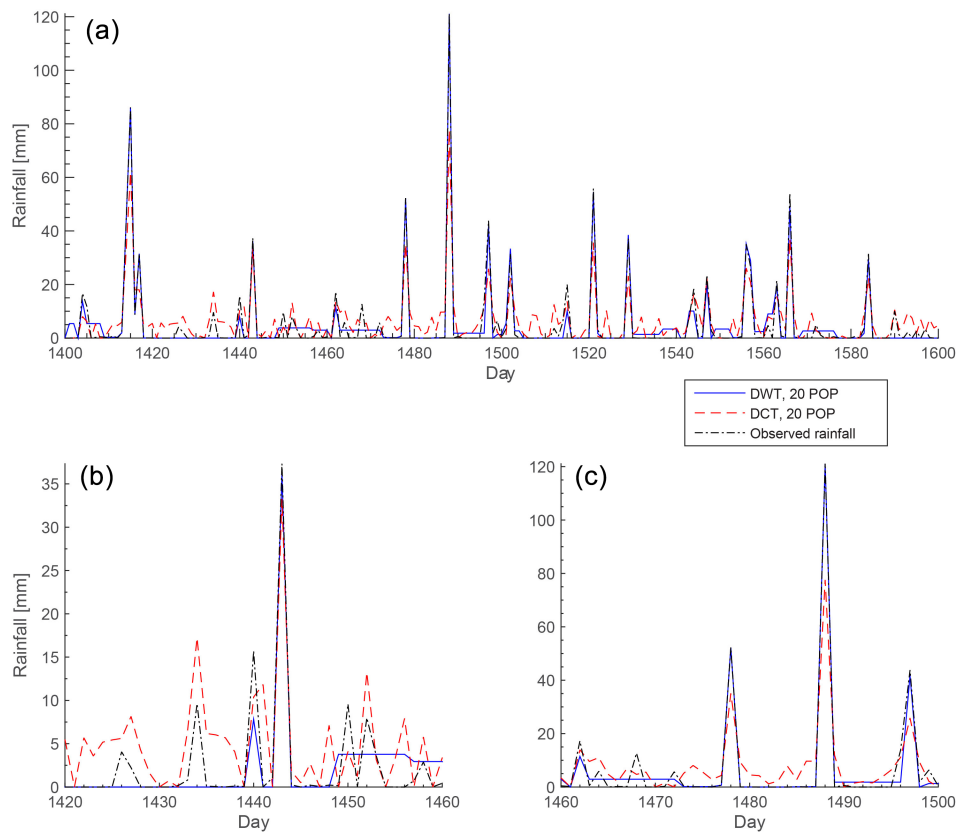
tion, the DWT will reconstruct temporal variances better than the DCT. The increased skewness for the reconstructed DWT signals compared to the observed input signals indicates that there is an increased reconstruction of low-magnitude rainfall events. On the contrary, the decreased normalized skewness for the reconstructed DCT signals indicates that a number of the low-magnitude rainfall events are tending to be reconstructed towards the mean. The kurtosis results shown in Fig. 6 demonstrate that, when compared to the observed input signal, events of extreme deviation explain more of the variance for the reconstructed DWT and less of the variance for the reconstructed DCT. Consequently, as the nature of the extreme deviations is a critical piece of information, the use of the DCT for model input data reduction for hydrologic studies that have short time steps and involving rainfall as an input is not recommended. It is also seen in Fig. 6 that the DCT is more likely to miss peak rainfall height information. Consequently, care needs to be taken when choosing a transform when peak height is critical. Further, the DCT should not be used for studies involving flood forecasting situations where the accuracy of peak height is critical.

Whilst it is important that rain gauges measure high-magnitude rainfall events with accuracy and precision, it is also important that low-magnitude rainfall events are recorded. Consequently, when evaluating the merits of the DCT and DWT to reconstruct rainfall it would be prudent to analyze the frequency in which each transform is either unable to reconstruct a rainfall event or erroneously constructs





**Figure 7.** Comparison of the reconstructed DCT and DWT signal for the Leaf River (Collins) catchment using 20 POP.



**Figure 8.** Panel (a) shows a time series comparison of the reconstructed DCT and DWT signals for the Leaf River (Collins) catchment using 20 POP for a period of 200 days. Panels (b) and (c) are smaller windows of the same time series during both low- and high-rainfall periods.

a rainfall event. Table 2 illustrates that, at times, both transforms will either fail to reconstruct a low-magnitude rainfall event or will erroneously construct a rainfall event when there was none observed in the original rainfall time series. In general, the DWT outperforms the DCT. The exception to this is at 10 POP. This is a result of the discrete nature of the DWT analysis function as opposed to the continuous analy-

sis function used in the DCT. As the POP decreases towards zero, both transforms miss more rainfall events.

Due to rapid increases in rainfall intensity, high-magnitude rainfall events tend to have high-frequency components. In Fig. 7, the smoothing of high-frequency, high-magnitude rainfall events by the DCT is made evident by the lower slope of the linear least squares fit for the DCT reconstruc-

**Table 2.** The mean and standard deviation (SD) for the number of missed rainfall events for the DWT and DCT using a different number of parameters.

POP	Number of missed rainfall events			
	DWT		DCT	
	Mean	SD	Mean	SD
40 %	239.004	117.317	587.934	155.375
30 %	398.005	138.793	645.495	159.378
20 %	581.591	145.769	696.288	168.524
10 %	852.340	168.590	748.075	184.910

tion of Leaf River observed rainfall data when compared to the DWT. This shows that the compressed DWT is able to retain more detail for high-magnitude rainfall events than the DCT. Using 20 POP, 730 DWT parameters are able to reconstruct observed rainfall with an RSR of 0.315, whereas 730 DCT parameters are able to reconstruct observed rainfall with an RSR of 0.540. Figures 7 and 8 show that the DWT often misses and sometimes smooths out low-magnitude rainfall events; the DCT, however, does reconstruct inaccurate rainfall at these times. Figure 8 also demonstrates that, at lower POPs, the DCT will smooth out and underestimate high-magnitude events whilst the DWT will maintain accuracy and precision.

## 7 Conclusions

Succinct descriptions of the DCT and DWT were provided to determine the suitability of each transform to be used as a tool for hydrologic model input data reduction. Due to their different construction, each transform provides different possibilities for use in model input data reduction. Since it is infeasible to estimate all transform parameters, the modeller could choose to estimate high- or low-frequency parameters of the DCT. This would result in minimal control of the temporal component being modified. Due to the multi-level decomposition of an input signal into high- and low-frequency parameters by the DWT, the modeller is able to specify the estimation of both time and frequency components. Hence, portions of the input data record can be targeted for estimation. The use of the DWT as a hydrologic model input data reduction technique allows the modeller more flexible options. A comparison of the DWTs' and DCTs' ability to reconstruct MOPEX rainfall data using standard simulation performance summary metrics, descriptive statistics and peak errors was then made, and it was found that the DWT is most efficient at preserving high-magnitude and transient rainfall events. Thus, it is recommended that the DWT be used as a model input data reduction technique for hydrologic studies that have short time steps and involve rainfall as an input. Considering that the bias for the reconstructed

## A. Wright et al.: Model input data reduction in hydrology

DWT rainfall signal is consistently lower than that of the reconstructed DCT signal and that the skewness, kurtosis and variance are also closest to the input rainfall signal, it is recommended that the DWT also be used as a model input data reduction technique for hydrologic studies that have long time steps with rainfall as an input.

*Data availability.* All data were obtained from publicly available data sets ([ftp://hydrology.nws.noaa.gov/pub/gcip/mopex/US\\_Data/](ftp://hydrology.nws.noaa.gov/pub/gcip/mopex/US_Data/); NWS, 2017).

*Author contributions.* AW conducted the experimental work, contributed towards the theory and wrote the paper. JW and DR assisted in the writing process. VP contributed towards the theory and assisted in the writing process.

*Acknowledgements.* The authors would like to extend their gratitude to Jasper Vrugt, Hamid Bazargan and the anonymous reviewers for their comments and recommendations. This work was supported by the Multi-modal Australian Sciences Imaging and Visualisation Environment (MASSIVE) (<http://www.massive.org.au>), a Monash University Engineering Research Living Allowance stipend and a top-up scholarship from the Bushfire & Natural Hazards Cooperative Research Centre. Valentijn Pauwels is funded by ARC grant FT130100545.

Edited by: Insa Neuweiler

Reviewed by: Hamid Bazargan and two anonymous referees

## References

- Ahmed, N., Natarajan, T., and Rao, K.: Discrete Cosine Transform, *IEEE T. Comput.*, C-23, 90–93, <https://doi.org/10.1109/TC.1974.223784>, 1974.
- Ajami, N., Duan, Q., and Sorooshian, S.: An integrated hydrologic Bayesian multimodel combination framework: Confronting input, parameter, and model structural uncertainty in hydrologic prediction, *Water Resour. Res.*, 43, W01403, <https://doi.org/10.1029/2005WR004745>, 2007.
- Beven, K.: A manifesto for the equifinality thesis, *J. Hydrol.*, 320, 18–36, <https://doi.org/10.1016/j.jhydrol.2005.07.007>, 2006.
- Blazkova, S. and Beven, K.: A limits of acceptability approach to model evaluation and uncertainty estimation in flood frequency estimation by continuous simulation: Skalka catchment, Czech Republic, *Water Resources Research*, 45, W00B16, <https://doi.org/10.1029/2007WR006726>, 2009.
- Bulygina, N. and Gupta, H.: Correcting the mathematical structure of a hydrological model via Bayesian data assimilation, *Water Resour. Res.*, 47, W05514, <https://doi.org/10.1029/2010WR009614>, 2011.
- Butts, M., Payne, J., Kristensen, M., and Madsen, H.: An evaluation of the impact of model structure on hydrological modelling uncertainty for streamflow simulation, *J. Hydrol.*, 298, 242–266, <https://doi.org/10.1016/j.jhydrol.2004.03.042>, 2004.

- Cloke, H. L. and Pappenberger, F.: Ensemble flood forecasting: A review, *J. Hydrol.*, 375, 613–626, 2009.
- Daubechies, I.: The Wavelet Transform, Time-Frequency Localization and Signal Analysis, *IEEE T. Inform. Theory*, 36, 961–1005, <https://doi.org/10.1109/18.57199>, 1990.
- De Vleeschouwer, N. and Pauwels, V. R. N.: Assessment of the indirect calibration of a rainfall-runoff model for ungauged catchments in Flanders, *Hydrol. Earth Syst. Sci.*, 17, 2001–2016, <https://doi.org/10.5194/hess-17-2001-2013>, 2013.
- Duan, Q., Sorooshian, S., and Gupta, V. K.: Optimal use of the SCE-UA global optimization method for calibrating watershed models, *J. Hydrol.*, 158, 265–284, 1994.
- Gan, T. Y. and Biftu, G. F.: Automatic calibration of conceptual rainfall-runoff models: Optimization algorithms, catchment conditions, and model structure, *Water Resour. Res.*, 32, 3513–3524, 1996.
- Kavetski, D., Kuczera, G., and Franks, S.: Bayesian analysis of input uncertainty in hydrological modeling: 2. Application, *Water Resour. Res.*, 42, W03408, <https://doi.org/10.1029/2005WR004376>, 2006a.
- Kavetski, D., Kuczera, G., and Franks, S.: Bayesian analysis of input uncertainty in hydrological modeling: 1. Theory, *Water Resour. Res.*, 42, W03407, <https://doi.org/10.1029/2005WR004368>, 2006b.
- Kuczera, G., Kavetski, D., Franks, S., and Thyer, M.: Towards a Bayesian total error analysis of conceptual rainfall-runoff models: Characterising model error using storm-dependent parameters, *J. Hydrol.*, 331, 161–177, <https://doi.org/10.1016/j.jhydrol.2006.05.010>, 2006.
- Kumar, P. and Foufoula-Georgiou, E.: Wavelet analysis for geophysical applications, *Rev. Geophys.*, 35, 385–412, 1997.
- Labat, D.: Recent advances in wavelet analyses: Part 1. A review of concepts, *J. Hydrol.*, 314, 275–288, 2005.
- Laloy, E. and Vrugt, J.: High-dimensional posterior exploration of hydrologic models using multiple-try DREAM (ZS) and high-performance computing, *Water Resour. Res.*, 48, W01526, <https://doi.org/10.1029/2011WR010608>, 2012.
- Mallat, S.: A Theory for Multiresolution Signal Decomposition: The Wavelet Representation, *IEEE T. Pattern Anal.*, 11, 674–693, <https://doi.org/10.1109/34.192463>, 1989.
- Mallat, S.: A Wavelet Tour of Signal Processing, third edition, Academic Press, Boston, USA, <https://doi.org/10.1016/B978-0-12-374370-1.50001-9>, 2009.
- McMillan, H., Jackson, B., Clark, M., Kavetski, D., and Woods, R.: Rainfall uncertainty in hydrological modelling: An evaluation of multiplicative error models, *J. Hydrol.*, 400, 83–94, 2011.
- Moriasi, D. N., Arnold, J. G., Van Liew, M. W., Bingner, R. L., Harmel, R. D., and Veith, T. L.: Model evaluation guidelines for systematic quantification of accuracy in watershed simulations, *T. ASABE*, 50, 885–900, 2007.
- Nalley, D., Adamowski, J., and Khalil, B.: Using discrete wavelet transforms to analyze trends in streamflow and precipitation in Quebec and Ontario (1954–2008), *J. Hydrol.*, 475, 204–228, <https://doi.org/10.1016/j.jhydrol.2012.09.049>, 2012.
- National Weather Service (NWS): Model Parameter Estimation Experiment (MOPEX), National Oceanic and Atmospheric Administration (NOAA), available at: [ftp://hydrology.nws.noaa.gov/pub/gcip/mopex/US\\_Data/](ftp://hydrology.nws.noaa.gov/pub/gcip/mopex/US_Data/), last access: 4 January 2017.
- Pauwels, V. and De Lannoy, G.: Multivariate calibration of a water and energy balance model in the spectral domain, *Water Resour. Res.*, 47, W07523, <https://doi.org/10.1029/2010WR010292>, 2011.
- Polikar, R.: The story of wavelets, Physics and modern topics in mechanical and electrical engineering, World Scientific and Engineering Society Press, Wisconsin, USA, 192–197, 1999.
- Renard, B., Kavetski, D., Leblois, E., Thyer, M., Kuczera, G., and Franks, S.: Toward a reliable decomposition of predictive uncertainty in hydrological modeling: Characterizing rainfall errors using conditional simulation, *Water Resour. Res.*, 47, W11516, <https://doi.org/10.1029/2011WR010643>, 2011.
- Robertson, D. E., Shrestha, D. L., and Wang, Q. J.: Post-processing rainfall forecasts from numerical weather prediction models for short-term streamflow forecasting, *Hydrol. Earth Syst. Sci.*, 17, 3587–3603, <https://doi.org/10.5194/hess-17-3587-2013>, 2013.
- Sadegh, M. and Vrugt, J.: Approximate Bayesian Computation using Markov Chain Monte Carlo simulation: DREAM(ABC), *Water Resour. Res.*, 50, 6767–6787, <https://doi.org/10.1002/2014WR015386>, 2014.
- Schaake, J., Cong, S., and Duan, Q.: The US mopex data set, in: Large Sample Basin Experiments for Hydrological Model Parameterization: Results of the Model Parameter Experiment – MOPEX, edited by: Andréassian, V., Hall, A., Chahinian, N., and Schaake, J., 9–28, IAHS Publ. no. 307, IAHS Press, Wallingford, UK, 2006.
- Schaeffli, B. and Zehe, E.: Hydrological model performance and parameter estimation in the wavelet-domain, *Hydrol. Earth Syst. Sci.*, 13, 1921–1936, <https://doi.org/10.5194/hess-13-1921-2009>, 2009.
- Shrestha, D., Robertson, D., Bennett, J., and Wang, Q.: Improving precipitation forecasts by generating ensembles through postprocessing, *Mon. Weather Rev.*, 143, 3642–3663, <https://doi.org/10.1175/MWR-D-14-00329.1>, 2015.
- Sivakumar, B.: Rainfall dynamics at different temporal scales: A chaotic perspective, *Hydrol. Earth Syst. Sci.*, 5, 645–652, <https://doi.org/10.5194/hess-5-645-2001>, 2001.
- Tang, Y., Reed, P., and Wagener, T.: How effective and efficient are multiobjective evolutionary algorithms at hydrologic model calibration?, *Hydrol. Earth Syst. Sci.*, 10, 289–307, <https://doi.org/10.5194/hess-10-289-2006>, 2006.
- Testik, F. Y. and Gebremichael, M. (Eds.): Rainfall: State of the Science, *Geoph. Monog. Series*, 191, 1–287, <https://doi.org/10.1029/gm191>, 2010.
- Thiemann, M., Trosset, M., Gupta, H., and Sorooshian, S.: Bayesian recursive parameter estimation for hydrologic models, *Water Resour. Res.*, 37, 2521–2535, 2001.
- Thyer, M., Kuczera, G., and Bates, B. C.: Probabilistic optimization for conceptual rainfall-runoff models: A comparison of the shuffled complex evolution and simulated annealing algorithms, *Water Resour. Res.*, 35, 767–773, 1999.
- Thyer, M., Renard, B., Kavetski, D., Kuczera, G., Franks, S., and Srikanthan, S.: Critical evaluation of parameter consistency and predictive uncertainty in hydrological modeling: A case study using Bayesian total error analysis, *Water Resour. Res.*, 45, W00B14, <https://doi.org/10.1029/2008WR006825>, 2009.
- Vrugt, J. A. and Robinson, B.: Treatment of uncertainty using ensemble methods: Comparison of sequential data assimilation and

- Bayesian model averaging, *Water Resour. Res.*, 43, W01411, <https://doi.org/10.1029/2005WR004838>, 2007.
- Vrugt, J. A. and Ter Braak, C. J. F.: DREAM<sub>(D)</sub>: an adaptive Markov Chain Monte Carlo simulation algorithm to solve discrete, noncontinuous, and combinatorial posterior parameter estimation problems, *Hydrol. Earth Syst. Sci.*, 15, 3701–3713, <https://doi.org/10.5194/hess-15-3701-2011>, 2011.
- Vrugt, J. A., Diks, C., Gupta, H., Bouten, W., and Verstraten, J.: Improved treatment of uncertainty in hydrologic modeling: Combining the strengths of global optimization and data assimilation, *Water Resour. Res.*, 41, 1–17, <https://doi.org/10.1029/2004WR003059>, 2005.
- Vrugt, J. A., Ter Braak, C., Clark, M., Hyman, J., and Robinson, B.: Treatment of input uncertainty in hydrologic modeling: Doing hydrology backward with Markov chain Monte Carlo simulation, *Water Resour. Res.*, 44, W00B09, <https://doi.org/10.1029/2007WR006720>, 2008.
- Vrugt, J. A., ter Braak, C., Gupta, H., and Robinson, B.: Equifinality of formal (DREAM) and informal (GLUE) Bayesian approaches in hydrologic modeling?, *Stoch. Env. Res. Risk A.*, 23, 1011–1026, <https://doi.org/10.1007/s00477-008-0274-y>, 2009a.
- Vrugt, J. A., ter Braak, C. J. F., Diks, C. G. H., Robinson, B. A., Hyman, J. M., and Higdon, D.: Accelerating Markov Chain Monte Carlo Simulation by Differential Evolution with Self-Adaptive Randomized Subspace Sampling, *Int. J. Nonlin. Sci. Num.*, 10, 273–290, 2009b.
- World Meteorological Organization: Guide to Meteorological Instruments and Methods of Observation, Geneva, Switzerland, 2014.
- Young, P.: Hypothetico-inductive data-based mechanistic modeling of hydrological systems, *Water Resour. Res.*, 49, 915–935, <https://doi.org/10.1002/wrcr.20068>, 2013.