

1 **Aggregation in environmental systems: Catchment mean** 2 **transit times and young water fractions under hydrologic** 3 **nonstationarity**

4
5 **J. W. Kirchner^{1,2}**

6 [1]{ETH Zürich, Zürich, Switzerland}

7 [2]{Swiss Federal Research Institute WSL, Birmensdorf, Switzerland}

8
9 Correspondence to: J. W. Kirchner (kirchner@ethz.ch)

10 11 **Abstract**

12 Methods for estimating mean transit times from chemical or isotopic tracers (such as Cl^- ,
13 $\delta^{18}\text{O}$, or $\delta^2\text{H}$) commonly assume that catchments are stationary (i.e., time-invariant) and
14 homogeneous. Real catchments are neither. In a companion paper, I showed that catchment
15 mean transit times estimated from seasonal tracer cycles are highly vulnerable to aggregation
16 error, exhibiting strong bias and large scatter in spatially heterogeneous catchments. I
17 proposed the young water fraction, which is virtually immune to aggregation error under
18 spatial heterogeneity, as a better measure of transit times. Here I extend this analysis by
19 exploring how nonstationarity affects mean transit times and young water fractions estimated
20 from seasonal tracer cycles, using benchmark tests based on a simple two-box model. The
21 model exhibits complex nonstationary behavior, with striking volatility in tracer
22 concentrations, young water fractions, and mean transit times, driven by rapid shifts in the
23 mixing ratios of fluxes from the upper and lower boxes. The transit-time distribution in
24 streamflow becomes increasingly skewed at higher discharges, with marked increases in the
25 young water fraction and decreases in the mean water age, reflecting the increased dominance
26 of the upper box at higher flows. This simple two-box model exhibits strong equifinality,
27 which can be partly resolved by simple parameter transformations. However, transit times are
28 primarily determined by residual storage, which cannot be constrained through hydrograph
29 calibration and must instead be estimated by tracer behavior.

1 Seasonal tracer cycles in the two-box model are very poor predictors of mean transit times,
2 with typical errors of several hundred percent. However, the same tracer cycles predict time-
3 averaged young water fractions (F_{yw}) within a few percent, even in model catchments that are
4 both nonstationary and spatially heterogeneous (although they may be biased by roughly 0.1-
5 0.2 at sites where strong precipitation seasonality is correlated with precipitation tracer
6 concentrations). Flow-weighted fits to the seasonal tracer cycles accurately predict the flow-
7 weighted average F_{yw} in streamflow, while unweighted fits to the seasonal tracer cycles
8 accurately predict the unweighted average F_{yw} . Young water fractions can also be estimated
9 separately for individual flow regimes, again with a precision of a few percent, allowing
10 direct determination of how shifts in hydraulic regime alter the fraction of water reaching the
11 stream by fast flowpaths. One can also estimate the chemical composition of idealized
12 "young water" and "old water" end-members, using relationships between young water
13 fractions and solute concentrations across different flow regimes. These results demonstrate
14 that mean transit times cannot be estimated reliably from seasonal tracer cycles, and that, by
15 contrast, the young water fraction is a robust and useful metric of transit times, even in
16 catchments that exhibit strong nonstationarity and heterogeneity.

17

18 Keywords: transit time, travel time, residence time, isotope tracers, residence time,
19 convolution, catchment hydrology, aggregation error, aggregation bias

20

21 **1 Introduction**

22 In a companion paper (Kirchner, 2015, hereafter referred to as Paper 1), I pointed out that
23 although catchments are pervasively heterogeneous, we often model them, and interpret
24 measurements from them, as if they were homogeneous. This makes our measurements and
25 models vulnerable to so-called "aggregation error", meaning that they yield inconsistent
26 results at different levels of aggregation. I illustrated this general problem with the specific
27 example of mean transit times (MTT's) estimated from seasonal tracer cycles in precipitation
28 and discharge. Using simple numerical experiments with synthetic data, I showed that these
29 MTT estimates will typically exhibit strong bias and large scatter when they are derived from
30 spatially heterogeneous catchments. Given that spatial heterogeneity is ubiquitous in real-
31 world catchments, these findings pose a fundamental challenge to the use of MTT's to
32 characterize catchment behavior.

1 In Paper 1 I also showed that seasonal tracer cycles in precipitation and streamflow can be
2 used to estimate the young water fraction F_{yw} , defined as the fraction of discharge that is
3 younger than a threshold age of approximately 2-3 months. I further showed that F_{yw}
4 estimates, unlike MTT estimates, are robust against extreme spatial heterogeneity. Thus
5 Paper 1 demonstrates the feasibility of determining the proportions of "young" and "old"
6 water (F_{yw} and $1-F_{yw}$, respectively) in spatially heterogeneous catchments.

7 But real-world catchments are not only heterogeneous. They are also nonstationary; their
8 travel-time distributions shift with changes in their flow regimes, due to shifts in the relative
9 water fluxes and flow speeds of different flowpaths (e.g., Kirchner et al., 2001; Tetzlaff et al.,
10 2007; Hrachowitz et al., 2010; Botter et al., 2010; Van der Velde et al., 2010; Birkel et al.,
11 2012; Heidbüchel et al., 2012; Peters et al., 2014). This nonstationarity is more than simply a
12 time-domain analogue to the heterogeneity problem explored in Paper 1, because variations in
13 flow regime may alter both the transit-time distributions of individual flowpaths and the
14 mixing ratios between them. Intuition suggests that catchment nonstationarity could play
15 havoc with estimates of MTT's, and perhaps also with estimates of the young water fraction.

16 This paper explores three central questions. First, does nonstationarity lead to aggregation
17 errors in MTT, and thus to bias or scatter in MTT estimates derived from seasonal tracer
18 cycles? Second, is the young water fraction F_{yw} also vulnerable to aggregation errors under
19 nonstationarity, or is it relatively immune, like it is to aggregation errors arising from spatial
20 heterogeneity? Third, can either MTT or F_{yw} be estimated reliably from seasonal tracer
21 cycles, in catchments that are both nonstationary and heterogeneous, as real catchments are?

22 In keeping with the spirit of the approach developed in Paper 1, here I explore the
23 consequences of catchment nonstationarity through simple thought experiments. These
24 thought experiments are based on a simple two-compartment conceptual model (Fig. 1). This
25 model greatly simplifies the complexities of real-world catchments, but it is sufficient to
26 illustrate the key issues at hand. It is not intended to simulate the behavior of a specific real-
27 world catchment, and thus its "goodness of fit" to any particular catchment time series is
28 unimportant. Instead, its purpose is to simulate how nonstationary dynamics may influence
29 tracer concentrations across wide ranges of catchment behavior, and thus to serve as a
30 numerical "test bed" for exploring how catchment nonstationarity affects our ability to infer
31 catchment transit times from tracer concentrations. One can of course construct more
32 complicated and (perhaps) realistic models, but that is not the point here. The point here is to

1 explore the consequences of catchment nonstationarity, in the context of one of the simplest
2 possible models which nonetheless exhibits a wide range of nonstationary behaviors.

3

4 **2 A simple conceptual model for exploring nonstationarity**

5 **2.1 Structure and basic equations**

6 The model catchment consists of two compartments, an upper box and a lower box (Fig. 1).
7 In typical conceptual models the upper box might represent soil water storage and the lower
8 box might represent groundwater, but for the present purposes it is unnecessary to assign the
9 two boxes to specific domains in the catchment. The upper box storage S_u is filled by
10 precipitation P , and drains at a leakage rate L that is a power function of storage; for
11 simplicity, evapotranspiration is ignored. Thus storage in the upper box evolves according to

$$12 \quad \frac{dS_u}{dt} = P - L = P - k_u S_u^{b_u} \quad , \quad (1)$$

13 where the coefficient k_u and the exponent b_u are parameters. A third parameter $0 < \eta < 1$
14 partitions the leakage L from the upper box into an amount ηL that flows directly to discharge
15 and an amount $(1-\eta)L$ that recharges the lower box. The lower box storage S_l is recharged by
16 leakage from the upper box and drains to streamflow at a discharge rate Q_l that is another
17 power function of storage,

$$18 \quad \frac{dS_l}{dt} = (1-\eta)L - Q_l = (1-\eta)L - k_l S_l^{b_l} \quad , \quad (2)$$

19 where the coefficient k_l and the exponent b_l are the final two parameters. The stream
20 discharge is the sum of the contributions from the upper and lower boxes, or

$$21 \quad Q_S = \eta L + Q_l \quad . \quad (3)$$

22 All storages are in mm of water equivalent depth, and all fluxes are in mm/day. The age
23 distribution in each box is explicitly tracked at daily resolution for the youngest 90 days, and
24 by accounting for the aggregate "age mass" (Bethke and Johnson, 2008) of each box's water
25 that is older than 90 days. The young water fraction F_{yw} is calculated as the fraction of water
26 in each box that is up to (and including) 69 days old; this threshold age equals 0.189 years,

1 which was shown in Paper 1 to be the theoretical young-water threshold age for seasonal
2 cycles in systems with exponential transit time distributions.

3 Discharge from both boxes is assumed to be non-age-selective, meaning that discharge is
4 taken proportionally from each part of the age distribution; thus the flow from each box will
5 have the same tracer concentration, the same young water fraction F_{yw} , and the same mean
6 age as the averages of those quantities in that box (at that moment in time). Tracer
7 concentrations and mean ages are tracked under the assumption that both boxes are each well-
8 mixed but also separate from one another, so their tracer concentrations and water ages will
9 differ. The tracer concentrations, young water fractions, and mean water ages in streamflow
10 are the flux-weighted averages of the contributions from the two boxes.

11 The model is solved on a daily time step, using a weighted combination of the partly implicit
12 trapezoidal method (for greater accuracy) and the fully implicit backward Euler method (for
13 guaranteed stability). Details of the solution scheme are outlined in Appendix A.

14 **2.2 Parameters and initialization**

15 The drainage coefficients k_u and k_l are problematic as model parameters, because their values
16 and dimensions are strongly dependent on the exponents b_u and b_l . Therefore I instead
17 parameterize the model drainage functions by the (dimensionless) exponents b_u and b_l and by
18 the (dimensional) "reference" storage values, $S_{u,ref}$ and $S_{l,ref}$. These reference values represent
19 the storage levels at which the drainage rates of each box will equal their long-term average
20 input rates. That is, $S_{u,ref}$ is the level of upper-box storage at which the leakage rate L equals
21 the long-term average input rate \bar{P} . Likewise, $S_{l,ref}$ is the level of lower-box storage at which
22 the discharge rate Q_l equals the average rate of recharge $(1 - \eta) \bar{L}$ (which, due to conservation
23 of mass in the upper box, also equals $(1 - \eta) \bar{P}$). The drainage function coefficients are
24 calculated from the reference storage values as follows:

$$25 \quad k_u S_{u,ref}^{b_u} = \bar{P} \quad , \quad k_u = \bar{P} S_{u,ref}^{-b_u} \quad (4)$$

$$26 \quad k_l S_{l,ref}^{b_l} = (1 - \eta) \bar{P} \quad , \quad k_l = (1 - \eta) \bar{P} S_{l,ref}^{-b_l} \quad . \quad (5)$$

27 Expressing k_u and k_l in this way is equivalent to writing the drainage equations for the two
28 boxes in dimensionless form, with the drainage rate expressed with reference to the long-term
29 input rate as follows:

$$1 \quad \frac{L}{\bar{P}} = \left(\frac{S_u}{S_{u,ref}} \right)^{b_u} \quad (6)$$

$$2 \quad \frac{Q_l}{(1-\eta)\bar{P}} = \left(\frac{S_l}{S_{l,ref}} \right)^{b_l} \quad (7)$$

3 One advantage of this approach is that, whereas the drainage coefficients k_u and k_l have no
 4 clear meaning and their numerical values and dimensions can vary wildly, the reference
 5 storage values are measured in mm of water equivalent depth, and their interpretation is
 6 straightforward. A further advantage of this approach is that it provides for varying degrees
 7 of residual storage without requiring any additional parameters to do so. Because $S_{u,ref}$ and
 8 $S_{l,ref}$ are the storage levels at which long-term mass balance is achieved, they represent the
 9 equilibria around which S_u and S_l will tend to fluctuate, with the range of those fluctuations
 10 largely determined by the variability in precipitation rates and by the stiffness of the drainage
 11 functions, as specified by the exponents b_u and b_l (see Sect. 3.2).

12 The storages are initialized at the reference values $S_{u,ref}$ and $S_{l,ref}$. The tracer concentrations
 13 are initialized at equilibrium (that is, at the volume-weighted mean of the precipitation tracer
 14 concentration). Likewise the mean ages in each box are initialized at their steady-state
 15 equilibrium values, $S_{u,ref} / \bar{P}$ in the upper box and $S_{u,ref} / \bar{P} + S_{l,ref} / [\bar{P}(1-\eta)]$ in the lower
 16 box. After a one-year spin-up period, I run the model for ten more years, with results for
 17 those ten years reported here.

18 **2.3 Parameter ranges and precipitation drivers**

19 Here I drive the model with three different real-world rainfall time series, representing a range
 20 of climatic regimes: a humid maritime climate with frequent rainfall and moderate seasonality
 21 (Plynlimon, Wales; Köppen climate zone Cfb), a Mediterranean climate marked by wet
 22 winters and very dry summers (Smith River, California, USA; Köppen climate zone Csb), and
 23 a humid temperate climate with very little seasonal variation in average rainfall (Broad River,
 24 Georgia, USA; Köppen climate zone Cfa). Figure 2 shows the contrasting frequency
 25 distributions and seasonalities of the three rainfall records. The Plynlimon rain gauge data
 26 were provided by the Centre for Ecology and Hydrology (UK), and the Smith River and
 27 Broad River precipitation data are reanalysis products from the MOPEX experiment [Duan,
 28 2006 #2193; ftp://hydrology.nws.noaa.gov/pub/gcip/mopex/US_Data/]. The use of these real-

1 world precipitation time series obviates the need to generate statistically realistic synthetic
2 precipitation to drive the model.

3 The model used here shares a similar overall structure with many other conceptual models
4 (e.g., Benettin et al., 2013), with several simplifications. But although the model used here is
5 typical in many respects, I will use it in an unusual way. Typically one calibrates a model to
6 reproduce the behavior of a real-world catchment, and then draws inferences about that
7 catchment from the parameters and behavior of the calibrated model. Here, however, the
8 model is not intended to represent any particular real-world system. Instead, the model itself
9 is the system under study, across wide ranges of parameter values, because the goal is to gain
10 insight into how nonstationarity affects general patterns of tracer behavior. Thus the fidelity
11 of the model in representing any particular catchment is not a central issue.

12 For the simulations shown here, the drainage exponents b_u and b_l are randomly chosen from
13 uniform distributions spanning the ranges of 1-20 and 1-50, respectively, the partitioning
14 coefficient η is randomly chosen from a uniform distribution ranging from 0.1 to 0.9, and the
15 reference storage levels $S_{u,ref}$ and $S_{l,ref}$ are randomly chosen from a uniform distribution of
16 logarithms spanning the ranges of 20-500 mm and 500-10,000 mm, respectively. These
17 parameter distributions are designed to encompass a wide range of possible behaviors,
18 including both strong and damped response to rainfall inputs, and small and large residual
19 storage. To illustrate the behavior of the model for one concrete case, I use a "reference"
20 parameter set with values taken from roughly the middle of each of these parameter
21 distributions ($b_u=10$, $b_l=20$, $\eta=0.5$, $S_{u,ref}=100$ mm, and $S_{l,ref}=2000$ mm). These parameter
22 values are not "better" than any others in any particular sense; they are simply a point of
23 reference (hence the name) for discussing the model's behavior.

24

25 **3 Results and discussion**

26 **3.1 Nonstationarity in the two-box model**

27 My main purpose is to use the simple two-box model to explore how catchment
28 nonstationarity affects our ability to infer water ages from tracer time series. I will take up
29 that issue beginning in Sect. 3.3 below. As background for that analysis, however, it is
30 helpful to first characterize the nonstationary behavior of the simple model system.

1 Figure 3 shows excerpts from the time series generated by the model with the Smith River
2 (Mediterranean climate) precipitation time series and the reference parameter set. One can
3 immediately see that the upper and lower boxes have markedly different mean ages (Fig. 3e),
4 young water fractions (Fig. 3d), and tracer concentrations (Fig. 3c), which also vary
5 differently through time. Tracer concentrations in the upper box (the orange line in Fig. 3c)
6 show a blocky, irregular pattern, remaining almost constant during periods of little rainfall,
7 and then changing rapidly when the box is episodically flushed by large precipitation events.
8 The lower box's tracer concentrations (the red line in Fig. 3c) are much more stable than the
9 upper box's, because its mean residence time is roughly 40 times longer ($S_{l,ref}$ is 20 times $S_{u,ref}$,
10 and with $\eta=0.5$, the flux through the lower box is only half of the flux through the upper box).
11 Because much more rain falls during the winters than the summers, the mean tracer
12 concentration in the lower box is closer to the winter concentrations than the summer
13 concentrations. During the wet winter season, rapid flushing keeps the young water fraction
14 near 100% in the upper box (the orange line in Fig. 3d), and can raise the young water
15 fraction to 30-40% in the lower box (the red line in Fig. 3d). Conversely, during the late
16 summer the young water fraction in the upper box temporarily dips to 50% or less, and the
17 young water fraction in the lower box declines to nearly zero. The small volume in the upper
18 box means that its water age (the orange line in Fig. 3e) is only a small fraction of a year. The
19 mean water age in the lower box (the red line in Fig. 3e) is much older and exhibits both
20 seasonal variation and inter-annual drift, reflecting year-to-year variations in total
21 precipitation. Thus the two components of this simple system have strongly contrasting
22 characteristics and behavior. These internal states of any real-world system would not be
23 observable, except as they are reflected in the volume and composition of streamflow.

24 In this regard, the most striking feature of Fig. 3 is the volatility of the tracer concentrations,
25 young water fractions, and mean transit times in discharge (the dark blue lines in Figs. 3c-e),
26 as the mixing ratio between the two boxes (Fig. 3b) shifts in response to precipitation events.
27 This mixing ratio is not a simple function of discharge (Fig. 4c); instead it is both hysteretic
28 and nonstationary, varying in response both to precipitation forcing and to the antecedent
29 moisture status of the two boxes (and thus to the prior history of precipitation). This
30 dependence on prior precipitation reflects the fact that the boxes typically retain their water
31 age and tracer signatures over time scales much longer than the timescale of hydraulic
32 response, because their residual storage is large compared to their dynamic storage (see Sect.
33 3.2). As a result, both the young water fraction and mean age of discharge and storage are

1 widely scattered functions of discharge (Figs. 4a and b). Likewise there is no simple
2 relationship between either the young water fraction or mean age in storage and the
3 corresponding quantities in discharge (Fig. 4d), although there is a strong overall bias toward
4 water in discharge being much younger than the average water in storage.

5 Even though drainage from each box is non-age-selective (that is, the young water fraction
6 and mean age in drainage from each box are identical to those in storage), this is emphatically
7 not true at the level of the two-box system, because the two boxes account for different
8 proportions of discharge than of storage. Furthermore, because the fractional contributions to
9 streamflow from the (younger, smaller) upper box and the (older, larger) lower box are highly
10 variable, the water age and young water fraction in discharge are not only strongly biased, but
11 also highly scattered, indicators of the same quantities in storage (Fig. 4d).

12 The aggregate long-term implications of these dynamics are evident in the marginal (time-
13 averaged) age distributions of storage and discharge (Fig. 5). From Fig. 5 it is immediately
14 obvious that the age distributions in discharge are strongly skewed toward young ages,
15 compared to the age distributions in storage, both for each box individually and for the
16 catchment as a whole. This skew toward young ages arises for two main reasons. First,
17 although drainage from each box is not age-selective, more outflow occurs during periods of
18 stronger precipitation forcing, and thus shorter residence times. Thus the average ages of the
19 outflow and the storage can differ greatly. Second, under high-flow conditions a larger
20 proportion of discharge is derived from the upper box (which has a relatively short transit
21 time), and at base flow more discharge is derived from the lower box (which has a larger
22 volume and a relatively long transit time). Thus the short-transit-time components of the
23 system dominate the discharge, while the long-transit-time components of the system
24 dominate the storage. As a result, the mean age in discharge will generally be much younger
25 than the mean age in whole-catchment storage, and likewise the young water fraction in
26 discharge will be much larger than the young water fraction in storage. Note that this is the
27 opposite of what one would expect from conceptual models like those of Botter (2012), in
28 which the mean water age in discharge either equals the mean age in storage (for well-mixed
29 systems), or is older than the mean age in storage (for piston-flow systems).

30 More generally, and more importantly, these results imply that estimates of water age in
31 streamflow cannot be translated straightforwardly into estimates of water age in storage.
32 Instead, they may underestimate the age of water in storage by large factors, although in the

1 particular example shown in Fig. 5, the difference is only about a factor of two. Three closely
2 related theoretical functions have recently been proposed to quantify the long-recognized
3 (Kreft and Zuber, 1978) disconnect between the age distributions in storage and in discharge.
4 These include the time-dependent StorAge Selection (SAS) function ω_Q of Botter et al.
5 (2011), the Storage Outflow Probability (STOP) functions of van der Velde et al. (2012), and
6 the rank StorAge Selection (rSAS) function of Harman (2015). While these functions are all
7 grounded in elaborate theoretical frameworks, it remains to be seen whether they can be
8 reliably estimated in practice using real-world data.

9 A further implication of the analysis above is that the marginal age distributions are not
10 exponential, even for individual boxes, and even though drainage from each box is not age-
11 selective. In steady state, non-age-selective drainage (i.e., the well-mixed assumption) would
12 yield an exponential distribution of ages in the upper box and in the short-time age
13 distribution in streamflow. However, when the system is not in steady state and we aggregate
14 its behavior over time, we are combining different age distributions from different moments
15 in time with different precipitation forcing. This creates an aggregation error in the time
16 domain, in the sense that the steady-state approximation will be a misleading guide to the
17 non-steady-state behavior of the system, *even on average*. That is, even over time scales
18 where inputs equal outputs and the long-term average fluxes are essentially constant – and
19 thus the steady-state approximation, on average, holds – the average behavior of the non-
20 steady-state system can differ significantly from the average behavior of an equivalent steady-
21 state system.

22 One can further explore these issues by examining the marginal (time-averaged) age
23 distributions for separate ranges of discharge (Fig. 6). Figure 6 shows that at higher
24 discharges, age distributions in streamflow are much more strongly skewed toward younger
25 ages, reflecting the increased dominance of the upper box at higher flows. For the upper half
26 of all discharges, the age distributions are more skewed than exponential; that is, they plot as
27 upward-curving lines in Fig. 6b. For the top 25% of discharges, they are approximately
28 power-law, plotting as nearly straight lines in Fig. 6c. The slopes of these lines are steeper
29 than 1, however, implying that the distributions must deviate from this trend at very short
30 ages; otherwise their integrals (i.e., their cumulative distributions) would become infinite. It
31 is important to note the mean ages quoted in Fig. 6a imply that the tails of the distributions all
32 extend far beyond the plot axes, which are truncated at 90 days. Note also that the

1 distributions shown in Fig. 6 have different shapes in different flow regimes, suggesting that
2 the model's high-flow behavior is not simply a re-scaled transform of its low-flow behavior.

3 **3.2 Residual storage and the disconnect between transit time and hydraulic** 4 **response time scales**

5 The model's complex, nonstationary water age and tracer dynamics arise from the disconnect
6 between the timescales of hydraulic response and catchment storage in each box, and from the
7 divergence in both these timescales between the two boxes. These contrasting timescales can
8 be estimated through simple scaling and perturbation analyses, as outlined in this section.

9 Total catchment storage consists of two components: the dynamic storage that is linked to
10 discharge fluctuations through storage-discharge relationships like Eqs. (6)-(7), plus the
11 residual or "passive" storage that remains when discharge has declined to very slow rates.

12 The range of dynamic storage exerts an important control on timescales of catchment
13 hydrologic response, while the much larger residual (or "passive") storage has little effect on
14 water fluxes but is an essential control on residence times (Kirchner, 2009; Birkel et al.,
15 2011).

16 In real-world catchments, sharply nonlinear storage-discharge relationships (Kirchner, 2009)
17 guarantee that dynamic storage will be small compared to residual storage. This behavior is
18 mirrored in the model, where if Eqs. (6)-(7) are strongly nonlinear (i.e., if the drainage
19 exponents b_u and b_l are much greater than 1), the volumes in the upper and lower boxes will
20 vary by only a small fraction of their reference storage values $S_{u,ref}$ and $S_{l,ref}$ (e.g., Fig. 3f).

21 They will remain relatively constant because, when the drainage exponents b_u and b_l are large,
22 the storage volumes cannot become much smaller than $S_{u,ref}$ and $S_{l,ref}$ without drainage rates
23 falling to near zero (thus stopping further decreases in storage), and conversely, the storage
24 volumes also cannot become much larger than $S_{u,ref}$ and $S_{l,ref}$ without drainage rates becoming
25 very high (thus stopping further increases in storage). Thus $S_{u,ref}$ and $S_{l,ref}$ will be a good
26 approximation to the residual storage volume, whenever the drainage exponents are much
27 greater than 1.

28 One can express this concept more quantitatively (though only approximately) using a simple
29 perturbation analysis. A first-order Taylor expansion of Eqs. (6) and (7) shows directly that
30 the fractional variability in drainage rates and storage are related by the drainage exponents in
31 the two boxes:

$$1 \quad \frac{\Delta L}{\bar{P}} \approx b_u \frac{\Delta S_u}{S_{u,ref}} \quad (8)$$

$$2 \quad \frac{\Delta Q_l}{(1-\eta)\bar{P}} \approx b_l \frac{\Delta S_l}{S_{l,ref}} \quad (9)$$

3 The variability in drainage rates from the upper and lower boxes, denoted as ΔL and ΔQ_l , will
4 be controlled by the temporal variability in precipitation; thus for a given precipitation
5 climatology, the dynamic variability in storage (denoted as ΔS_u and ΔS_l) will scale according
6 to the ratios $S_{u,ref}/b_u$ and $S_{l,ref}/b_l$. For example, when the model is driven by Smith River
7 precipitation and uses the reference parameters (Fig. 3), the variability in discharge from the
8 lower box, as measured by its standard deviation, is 3.7 mm d^{-1} , nearly equal to the average
9 lower box discharge of 3.8 mm d^{-1} . Because the reference value of b_l is 20, Eq. (9) implies
10 that the standard deviation of lower box storage should be approximately $1/20^{\text{th}}$ of the
11 reference storage $S_{l,ref}$, or roughly 100 mm. Consistent with this estimate, the actual standard
12 deviation of S_l is 84 mm or about 4% of the total. Figure 3f shows that at least 90% of $S_{l,ref}$ is
13 residual storage that never drains during the 10-year simulation, roughly consistent with the
14 perturbation analysis.

15 The perturbation analysis also yields estimates for the time scale of hydraulic response (which
16 controls how "flashy" the discharge will be), through a rearrangement of Eqs. (8) and (9) as
17 follows:

$$18 \quad \frac{\Delta S_u}{\Delta L} \approx \frac{S_{u,ref}}{b_u \bar{P}} \quad (\text{hydraulic response time scale, upper box}) \quad (10)$$

$$19 \quad \frac{\Delta S_l}{\Delta Q_l} \approx \frac{S_{l,ref}}{b_l (1-\eta) \bar{P}} \quad (\text{hydraulic response time scale, lower box}) \quad (11)$$

20 Again using the reference parameter values and Smith River precipitation (for which \bar{P} is
21 roughly 7.6 mm d^{-1}), Eqs. (10) and (11) imply a hydraulic response time of roughly 1.3 days
22 (for $b_u=10$) in the upper box and roughly 26 days (for $b_l=20$) in the lower box. These time
23 scales are factors b_u and b_l smaller than the steady-state mean transit times, which are
24 determined by the ratios between the volumes and water fluxes,

$$25 \quad \frac{S_{u,ref}}{\bar{P}} \quad (\text{steady-state mean transit time, upper box}) \quad (12)$$

$$1 \quad \frac{S_{l,ref}}{(1-\eta)\bar{P}} \quad (\text{steady-state mean transit time, lower box}) \quad (13)$$

2 From Eqs. (12)-(13) one can also directly estimate the steady-state mean travel time in the
 3 combined discharge as the weighted average of streamflow derived directly from the upper
 4 box, and water that flows through the upper and lower boxes in series,

$$5 \quad \eta \frac{S_{u,ref}}{\bar{P}} + (1-\eta) \left(\frac{S_{u,ref}}{\bar{P}} + \frac{S_{l,ref}}{(1-\eta)\bar{P}} \right) = \frac{S_{u,ref} + S_{l,ref}}{\bar{P}} \quad (14)$$

6 which is the expected result for any system at steady state: regardless of its internal
 7 configuration, the mean transit time in any steady-state system will equal the ratio between its
 8 storage volume and its throughput rate. For the reference parameter set and Smith River
 9 precipitation, Eq. (14) becomes $(100 \text{ mm} + 2000 \text{ mm})/7.6 \text{ mm d}^{-1}$, or roughly 0.76 years, in
 10 good agreement with the whole-catchment mean transit time of 0.74 years determined from
 11 age tracking (see Fig. 5d). Note, however, that the *distribution* of these transit times will be
 12 markedly different from the exponential distribution that would be expected in steady state.
 13 This makes estimating mean transit times from tracer fluctuations difficult, as shown below in
 14 Sect. 3.3.

15 Equations (12) and (13) imply that the mean transit times in the upper and lower boxes should
 16 be roughly 13 days (or 0.036 yr) and 529 days (or 1.45 yr), respectively, in good agreement
 17 with the mean transit times of 0.03 and 1.44 years determined from age tracking (Fig. 5d).
 18 However, Eqs. (10)-(11) imply that these transit times will differ by factors of 10 and 20 (the
 19 values of b_u and b_l , respectively) from the hydraulic response timescales that regulate
 20 catchment runoff response. The disconnect between hydraulic response times and mean
 21 transit times is the counterpart, in lumped conceptual models, to the disconnect between the
 22 velocity of water transport and the celerity of hydraulic head propagation in more realistic,
 23 physically extended systems (Beven, 1982; Kirchner et al., 2000; McDonnell and Beven,
 24 2014). This contrast between hydraulic response times and mean transit times (or dynamic
 25 and total storage, or celerity and velocity) is a simple explanation for the apparent paradox of
 26 prompt discharge of old water during storm events (Kirchner, 2003).

3.3 Inferring MTT and F_{yw} from seasonal tracer cycles in nonstationary catchments

The analysis above shows that the simple two-box model gives hydrograph and tracer behavior that is complex and nonstationary (Figs. 3-6). Furthermore, even this simple five-parameter model exhibits strong equifinality (Appendix B). Much of this equifinality can be alleviated (compare Figs. B1 and B2) through parameter transformations based on the perturbation analysis outlined above. However, because the timescales of catchment storage and hydraulic response are controlled by different combinations of parameters, parameter calibration to the hydrograph cannot constrain the storage volumes or streamwater age (Figs. B2-B3). These model results demonstrate general principles that have been recognized for years: a) the hydrograph responds to, and thus can help to constrain, dynamic storage but not passive storage, and b) because passive storage is often large, timescales of hydrologic response and catchment water storage are decoupled from one another, such that water ages cannot be inferred from hydrograph dynamics. Thus for understanding how catchments store and mix water, tracer data are essential.

But how should these tracer data be used? One approach is to explicitly include tracers in a catchment model, and calibrate that model against both the hydrograph and the tracer chemograph (e.g., Birkel et al., 2011; Benettin et al., 2013; Hrachowitz et al., 2013). The usefulness of that approach depends on whether the model parameters can be constrained and, more importantly, whether the model structure adequately characterizes the system under study (which is usually unknown, and possibly unknowable). Except in multi-model studies, it will be unclear how much the conclusions depend on the particular model that was used, and the particular way that it was fitted to the data. Furthermore, adequate tracer data for calibrating such models are rare, particularly because dynamic models require input data with no gaps. The mismatch between model complexity and data availability means that in some cases, all the data are used for calibration and validation must be skipped, leaving the reproducibility of the model results unclear (e.g., Benettin et al., 2015).

For all of these reasons, there will be an ongoing need for methods of inferring water ages that have modest data requirements and that are not dependent on specific model structures and parameters. Sine-wave fitting of seasonal tracer cycles, for example, is not based on a particular mechanistic model, but instead is based on a broader conceptual framework in which stream output is some convolution of previous precipitation inputs. That premise is of

1 course open to question, but nevertheless seasonal tracer cycles (of, e.g., ^{18}O , ^2H , and Cl^-)
2 have been widely used to estimate mean catchment transit times (see McGuire and
3 McDonnell, 2006, and references therein), largely because this particular method has modest
4 data requirements. In particular, it does not need unbroken records of either precipitation
5 inputs or streamflow outputs.

6 As detailed more fully in Paper 1, the seasonal tracer cycle method is based on the principle
7 that when one convolves a sinusoidal tracer input with a transit time distribution (TTD), one
8 obtains a sinusoidal output that is damped and phase-lagged by an amount that depends on the
9 shape of the TTD and also on its scale, as expressed, for example, by its mean transit time
10 (MTT). Conventionally one assumes an exponential TTD, which is the steady-state solution
11 for a well-mixed reservoir. More generally, one might assume that transit times are gamma-
12 distributed, recognizing that the exponential distribution is a special case of the gamma
13 distribution (with the shape factor α equal to 1). A sinusoidal tracer cycle that has been
14 convolved with a gamma TTD will be damped and phase-lagged as described in Eqs. (8) and
15 (9) of Paper 1. These equations can then be inverted to infer the shape and scale of the TTD
16 from the seasonal tracer cycles in precipitation and streamflow.

17 The procedure is as follows. One first measures the amplitudes and phases of the seasonal
18 tracer cycles in precipitation and streamflow using Eqs. (4)-(6) of Paper 1. If one assumes an
19 exponential TTD, one can estimate the MTT directly from the amplitude ratio A_S/A_P in
20 streamflow and precipitation using Eq. (10) of Paper 1 with $\alpha=1$. Where I plot results from
21 this procedure (i.e., Fig. 7) the corresponding axis will say "MTT inferred from A_S/A_P ". This
22 is the approach that is conventionally used in the literature. Alternatively, as I showed in
23 Sect. 4.4 of Paper 1, one can use the tracer cycle amplitude ratio A_S/A_P and phase shift $\phi_S - \phi_P$
24 to jointly estimate the shape factor α and the MTT (assuming the TTD is gamma-distributed,
25 which is less restrictive than assuming that it is exponential). To do this one estimates the
26 shape factor α from A_S/A_P and $\phi_S - \phi_P$ using Eq. 11 from Paper 1, and then estimates the scale
27 factor β using Eq. 10 from Paper 1; the MTT is α times β . MTT's estimated by this procedure
28 are shown in Figs. 10-12 as "MTT inferred from A_S/A_P and $\phi_S - \phi_P$ ".

29 Paper 1 shows that both of these MTT measures are extremely vulnerable to aggregation bias
30 in spatially heterogeneous catchments. Therefore Paper 1 proposes an alternative measure of
31 travel times, the young water fraction F_{yw} , which is designed to be much less sensitive than
32 MTT to aggregation artifacts. F_{yw} is the fraction of streamflow that is younger than a

1 specified threshold age. For a seasonal cycle (i.e., with a period of 1 year) and reasonable
2 range of TTD shapes, the threshold age varies between about 0.15 and 0.25 years, or
3 equivalently ~2-3 months (see Eq. 14 and Fig. 10 in Paper 1). As described in Sect. 2, in the
4 model simulations the "true" F_{yw} is defined by a threshold age of 0.189 years (69 days), which
5 equals the threshold age for seasonal cycles convolved with an exponential TTD.

6 One can use seasonal tracer cycles to infer the young water fraction following either of two
7 strategies. As shown in Sect. 4.1 of Paper 1, in many situations F_{yw} is approximately equal to
8 the amplitude ratio A_S/A_P itself (indeed, it was designed to have this property). In figures
9 where the amplitude ratio A_S/A_P is used as an estimate of F_{yw} (e.g., Fig. 7), the axis says
10 simply " F_{yw} inferred from A_S/A_P ". Alternatively, one can use both the amplitude ratio A_S/A_P
11 and phase shift $\phi_S-\phi_P$ to estimate F_{yw} , as explained in Sect. 4.4 of Paper 1. First, one estimates
12 the shape factor α from A_S/A_P and $\phi_S-\phi_P$ using Paper 1's Eq. (11). One then determines the
13 threshold age τ_{yw} from α using Paper 1's Eq. (14), and the scale factor β from α and A_S/A_P
14 using Paper 1's Eq. (10). Lastly, one estimates F_{yw} as lower incomplete gamma function
15 $\Gamma(\tau_{yw}, \alpha, \beta)$ (Eq. 13 of Paper 1). Where I have followed this more complex procedure (e.g.,
16 Figs. 9-12), the figure axes say " F_{yw} inferred from A_S/A_P and $\phi_S-\phi_P$ ". All of these F_{yw} 's and
17 MTT's are intended as temporal averages, reflecting whatever conditions (e.g., precipitation
18 climatologies or flow regimes) have shaped the seasonal cycles that are used to estimate them.

19 These methods for inferring the young water fraction F_{yw} are derived from the properties of
20 gamma TTD's. However, as I showed in Sects. 4.2-4.3 of Paper 1, these methods reliably
21 estimate F_{yw} for very wide ranges of catchment TTD's (beyond the already broad family of
22 gamma distributions), at least in catchments that are spatially heterogeneous but time-
23 invariant. Here I explore whether these methods are also reliable in nonstationary catchments
24 (and, in Sect. 3.5 below, in catchments that are both nonstationary and spatially
25 heterogeneous).

26 Figure 7 shows the true young water fractions F_{yw} and mean transit times (MTT's) in
27 discharge from the two-box model, compared to estimates of F_{yw} and MTT inferred from the
28 model's seasonal tracer cycles. As Figs. 7a-c show, the amplitude ratios A_S/A_P of seasonal
29 tracer cycles reliably estimate the true young water fractions in the model streamflow, across
30 1000 random parameter sets encompassing a very wide range of nonstationary catchment
31 behavior. The slight underestimation bias in Figs. 7a-c is reduced when both amplitude and
32 phase information are used to estimate F_{yw} (Figs. 7d-f). Under strongly seasonal precipitation

1 forcing (Smith River; right panels in Fig. 7), the seasonal tracer cycles underestimate F_{yw} by
2 roughly 0.1 to 0.2, although the predicted and observed values of F_{yw} remain strongly
3 correlated. For the other two precipitation drivers (Broad River and Plynlimon), the predicted
4 and observed values of F_{yw} correspond almost exactly. Thus Fig. 7 shows that the young
5 water fraction is relatively insensitive to aggregation error under nonstationarity, mirroring its
6 robustness against spatial heterogeneity (as shown in Paper 1). By contrast, estimates of MTT
7 are strongly biased and widely scattered, even on logarithmic axes (lower panels, Fig. 7).

8 One additional complication in nonstationary situations, compared to the time-invariant
9 examples explored in Paper 1, is that the young water fraction F_{yw} and mean transit time MTT
10 can be expressed either as simple averages over time (representing the F_{yw} or MTT of an
11 average *day* of streamflow), or as flow weighted averages (representing the F_{yw} or MTT of an
12 average *liter* of streamflow). These quantities will not be equivalent, since higher flows will
13 typically have higher F_{yw} 's and shorter MTT's (Figs. 3 and 4). Likewise one can expect that
14 amplitudes of flow-weighted and un-weighted fits to the seasonal tracer cycles will be
15 different. As the light blue points in Fig. 7 show, amplitude ratios of flow-weighted fits to the
16 seasonal tracer cycles accurately predict the flow-weighted F_{yw} in streamflow; likewise, as the
17 dark blue points show, the amplitude ratios of un-weighted fits accurately predict the un-
18 weighted F_{yw} in streamflow. The flow-weighted fits to the seasonal tracer cycles were
19 calculated by weighted least squares, with weights proportional to streamflow or precipitation
20 volume. (In real-world applications, a robust fitting technique like Iteratively Reweighted
21 Least Squares (IRLS) can be used to limit the influence of outliers. An R script for
22 performing volume-weighted IRLS is available from the author.)

23 The underestimation bias in F_{yw} observed under the Smith River precipitation forcing may
24 arise because the assumed tracer cycle is correlated with the strong seasonality in
25 precipitation, such that tracer concentrations peak during the summer, when almost no rain
26 falls. Thus the effective variability of tracer inputs to the catchment is less than one would
27 infer from a sinusoidal fit to the precipitation tracer concentrations (and volume-weighting the
28 fit does not help, because in these synthetic precipitation data the fit is exact so there are no
29 residuals on which the weighting can have any effect). Because the tracer concentration
30 amplitude overestimates the effective variability in tracer concentrations reaching the
31 catchment, the tracer damping in the catchment is overestimated and thus the F_{yw} is
32 underestimated. This underestimation bias disappears if one shifts the phase of the assumed

1 precipitation tracer concentrations so that they peak in the spring or fall, and thus are
2 uncorrelated with the seasonality in precipitation volumes. I have not done so here, however,
3 because stable isotope ratios in precipitation typically peak in the mid-summer at latitudes
4 poleward of $\sim 35^\circ$ (Feng et al., 2009), where most catchment studies have been conducted.
5 Thus Fig. 7 suggests the potential for bias in F_{yw} estimates at sites where isotope cycles are
6 correlated with very strong precipitation seasonality. However, even under the strongly
7 seasonal Smith River precipitation forcing, the bias in inferred F_{yw} values is small compared
8 to the a priori uncertainty in F_{yw} (which is of order 1), and small compared to the bias in
9 inferred MTT's (which is large even on logarithmic axes).

10 Figures 7g-i compare the MTT in streamflow with estimates of MTT as they are
11 conventionally calculated, that is, from the seasonal tracer cycle amplitude assuming an
12 exponential TTD. These figures show that these conventional estimates are subject to a
13 strong underestimation bias, which can exceed an order of magnitude. Some of the MTT
14 estimates do fall close to the 1:1 line, but these are mostly cases in which the partition
15 coefficient η is very small, such that nearly all drainage from the upper box is routed through
16 the lower box, thus transforming the two-box, nonstationary model into a nearly one-box,
17 nearly stationary model. The strong aggregation bias in MTT under catchment
18 nonstationarity shown in Figs. 7g-i mirrors the similarly strong bias under spatial
19 heterogeneity that was demonstrated in Paper 1.

20 The implication of Figs. 7g-i (and of Paper 1) is that many of the MTT values in the literature
21 are likely to be underestimated by large factors, and thus that real-world catchment MTT's are
22 likely to be much longer than we thought. This observation raises the question: where is all
23 that water being stored? In steady state, the storage volume must equal the discharge
24 multiplied by the MTT (see Sect. 3.2). Thus if we have been underestimating MTT's by large
25 factors, then we have also been underestimating catchment storage volumes by similar
26 multiples. Where is the storage volume that can accommodate all this water?

27 One possible answer is that in a non-steady-state system, the MTT decreases with increasing
28 discharge (e.g., Fig. 4b), and the storage volume equals the discharge multiplied by the
29 *volume-weighted* MTT rather than the *time-averaged* MTT. Because the volume-weighted
30 MTT is less (potentially much less) than the time-averaged MTT (see also Peters et al., 2014),
31 the implied storage volume is correspondingly smaller. Furthermore, many MTT studies in
32 the literature have been based on tracer sampling that excludes high flows, such that they infer

1 the mean age of baseflow rather than of the average discharge (McGuire and McDonnell,
2 2006). To the extent that mean baseflow discharges are lower than mean total discharges, the
3 stored volume of baseflow water will be less than what one might overestimate by
4 multiplying the mean *total* discharge by the mean *baseflow* age. Beyond these general
5 considerations, however, it makes little sense to draw precise inferences based on MTT
6 estimates that are likely to be strongly biased and widely scattered (as shown here, and also in
7 Paper 1).

8 It is important to recognize that the predicted F_{yw} values are really predictions, unlike many
9 "predictions" from calibrated models. The horizontal axes in Fig. 7 are calculated solely from
10 the age-tracking within the model, with no information about the tracer concentrations.
11 Likewise the vertical axes in Fig. 7 are calculated from the modeled tracer cycles alone,
12 without any information about the model that generated them, and in particular without any
13 information about the modeled age of streamflow. Thus Fig. 7 gives some basis for
14 confidence that estimates of F_{yw} will also be reliable in real-world catchments, where the true
15 "model" can never be known.

16 **3.4 Young water fractions in discrete flow regimes**

17 Figures 3 and 4 show that high-flow periods are characterized by shorter mean transit times
18 and higher young water fractions, reflecting the increased dominance of drainage from the
19 upper box with its younger water ages. Although instantaneous transit time distributions
20 (TTD's) can be highly variable, and thus instantaneous mean transit times and young water
21 fractions can exhibit scattered relationships with discharge (Fig. 4), the marginal (time-
22 averaged) TTD's in Fig. 6 clearly show systematically stronger skew toward younger water
23 ages in higher ranges of streamflow. Thus, as Fig. 6 shows, the TTD varies in shape, not just
24 in scale, between different flow regimes.

25 This observation leads naturally to the question of whether these variations in TTD's are also
26 reflected in streamflow tracer concentrations, and whether those tracer signatures can be used
27 to draw inferences about the TTD's that characterize individual flow regimes. Figure 3 shows
28 that high-flow periods typically exhibit wider variations in tracer concentrations, reflecting
29 greater contributions from the upper box, which has shorter residence times and thus more
30 labile tracer concentrations than the lower box does. To test how systematic these variations
31 in concentrations are, I ran the model with the reference parameter set and Plynlimon

1 (temperate maritime) precipitation forcing, and separated the resulting time series into six
2 discharge ranges. Figure 8 shows these six discharge ranges and the corresponding tracer
3 concentrations in dark blue, superimposed on the entire discharge and concentration time
4 series in light gray. As Fig. 8 shows, seasonal tracer cycles at higher flows are systematically
5 less damped and phase-shifted (relative to the tracer cycle in precipitation, shown by the
6 dotted gray line), implying shorter MTT's and larger young water fractions.

7 To test whether these changes in the seasonal tracer cycles are quantitatively consistent with
8 the shifts in water age across the six flow regimes, I fitted sinusoids separately to the tracer
9 concentrations in each individual discharge range (Fig. 8). I compared these with a single
10 sinusoid fitted to the entire precipitation tracer time series (because it is not possible to assign
11 discrete precipitation events to individual discharge ranges). From the resulting amplitude
12 ratios and phase shifts for each discharge range, I then estimated F_{yw} and MTT using the
13 methods outlined in Sect. 3.3. Figure 9 presents the results of this thought experiment,
14 showing that the time-averaged (but flow-specific) young water fraction F_{yw} in each discharge
15 range is accurately predicted by the damping and phase shift of the corresponding seasonal
16 tracer cycle.

17 To test whether this result is general, I repeated this thought experiment for 200 random
18 parameter sets and all three precipitation drivers. The results are shown in Fig. 10, with each
19 discharge range plotted in a different color. The colors overlap because the discharge ranges,
20 F_{yw} 's and MTT's all vary substantially from one parameter set to the next. The amplitudes and
21 phase shifts of the seasonal tracer cycles predict the time-averaged young water fractions F_{yw}
22 in each discharge range with reasonable accuracy (upper panels, Fig. 10). Somewhat
23 surprisingly, the F_{yw} underestimation bias seen in Figs. 7c and 7f under the highly seasonal
24 Smith River precipitation forcing does not arise in the predicted F_{yw} values for the separate
25 discharge ranges (Fig. 10c). In contrast to the generally close correspondence between the
26 predicted and observed F_{yw} values, predicted MTT's are very widely scattered for all
27 discharge ranges and all precipitation forcings (lower panels, Fig. 10).

28 **3.5 Combined effects of nonstationarity and spatial heterogeneity**

29 Paper 1 explored whether mean travel times and young water fractions can be reliably
30 inferred from tracer dynamics in spatially heterogeneous (but stationary) catchments,
31 composed of diverse subcatchments with different (but time-invariant) TTD's. The sections

1 above have presented a similar analysis for nonstationary (but spatially homogeneous)
2 catchments. However, real-world catchments are not *either* heterogeneous *or* nonstationary;
3 instead they are *both* heterogeneous *and* nonstationary. That is, their subcatchments each
4 exhibit nonstationary dynamics that may vary greatly from one to the next. To explore the
5 combined effects of nonstationarity and spatial heterogeneity, I merged the approach
6 developed in Paper 1 with the model developed in Sect. 2 above.

7 As illustrated in Fig. 11, I ran eight copies of the nonstationary model developed in Sect. 2,
8 representing eight different tributaries, each with a different, randomly chosen parameter set.
9 I chose the number eight to provide a reasonable degree of complexity and heterogeneity
10 while preserving a reasonable degree of computational efficiency. I supplied the same
11 precipitation forcing (Fig. 11a) to all eight models (Fig. 11b) to simulate the behavior of the
12 eight hypothetical tributary streams (Fig. 11c). I then simulated the merging of these streams
13 by averaging their discharges, and taking volume-weighted averages of their tracer
14 concentrations, young water fractions, and water ages (Fig. 11d). Because the instantaneous
15 flows from the eight tributaries vary differently through time, their mixing ratios also
16 fluctuate. The individual random parameter sets create a wide range of model structures at
17 the whole-catchment level, since the eight parallel subcatchments in Fig. 11 jointly comprise
18 a 16-box, 40-parameter model incorporating wide ranges of large and small reservoirs with
19 varying degrees of nonlinearity.

20 In any spatially heterogeneous catchment (which is to say, any real-world catchment), one
21 will typically only have observations from the merged whole-catchment streamflow (i.e., the
22 blue time series in Fig. 11d). One will typically have no information about the behavior of
23 the individual tributaries (i.e., the colored time series in Fig. 11c), and if one did, then those
24 tributaries would themselves have their own spatially heterogeneous tributary streams or
25 flowpaths, and so on. Thus the heterogeneity of any real-world catchment will remain poorly
26 quantified (and possibly even unrecognized), and rigorously reductionist attempts to fully
27 characterize such complex multiscale heterogeneity would be impractical.

28 Thus we face the problem: how much can we infer from the behavior of the merged whole-
29 catchment streamflow, given that it originates from processes that are heterogeneous and
30 nonstationary (to a degree that is unknown and unknowable)? Figure 12 explores this general
31 question in the specific context of young water fractions and mean travel times, presenting
32 results from 200 iterations of the heterogeneous nonstationary model shown in Figure 11 with

1 all three precipitation drivers. In Fig. 12 the merged streamflow is separated into discrete
2 flow regimes, following the approach outlined in Sect. 3.4. As Fig. 12 shows, F_{yw} values
3 inferred from the tracer cycles in each discharge range accurately predict the true fraction of
4 young water in that discharge range, as determined from age tracking.

5 Figure 12 is analogous to Fig. xxf10, with the difference that Fig. 10 shows model runs for
6 individual random parameter sets, whereas Fig. 12 shows results from eight runs merged
7 together. Merging the model outputs will tend to average out the idiosyncrasies of the
8 individual parameter sets, which is why the clusters of points in Fig. 12 are more compact
9 than the corresponding point clouds in Fig. 10. As a result, the individual discharge ranges
10 overlap less in Fig. 12 than in Fig. 10. The compact scatterplots shown in Fig. 12 show only
11 small deviations from the 1:1 line for estimates of the young water fraction F_{yw} . By contrast,
12 estimates of mean transit times in Fig. 12 exhibit substantial bias and scatter (note the
13 logarithmic axes in Fig. 12d-f).

14 **3.6 Hydrological and hydrochemical implications of young water fractions**

15 The results reported above, together with the results reported in Paper 1, show that unlike
16 mean transit times, young water fractions can be estimated reliably from seasonal tracer
17 cycles in catchments that are spatially heterogeneous, nonstationary, or both. These findings
18 then raise the obvious question: we can measure young water fractions reliably, but what are
19 they good for? One answer is that young water fractions can be considered as a catchment
20 characteristic, analogous (but far from equivalent) to MTT. In theory MTT should be
21 particularly useful as a catchment descriptor, because the MTT times the mean annual
22 discharge yields the total catchment storage. But because estimates of MTT will often be
23 substantially in error, estimates of catchment storage derived from MTT are likely to be
24 equally unreliable. If the shape of the transit time distribution (TTD) were known, of course,
25 there would be a clear functional relationship between MTT and F_{yw} , and one could be
26 calculated from the other. But if the shapes of the TTD were known, estimating the MTT
27 itself would also be easy; the problem in estimating the MTT is the fact that the TTD's shape
28 – particularly the length of its tail – is poorly constrained by tracer data. This is why F_{yw} can
29 be estimated much more reliably than MTT. F_{yw} , like the amplitude of the seasonal tracer
30 cycle, depends on the relative proportions of younger and older water, but is insensitive to
31 how old the "older" water is. MTT depends critically on the age of the older water, which

1 cannot be reliably determined because it has almost no effect on the seasonal tracer cycle (or
2 on more elaborate convolution analyses: see Seeger and Weiler, 2014).

3 Because the young water fraction is indifferent to the age of the older water, it cannot be used
4 to estimate residual storage. What F_{yw} estimates, instead, is the fraction of water reaching the
5 stream by relatively fast (less than ~2-3 month) flowpaths. In the context of the present
6 model, this is reflected in the correlation between F_{yw} and the partitioning parameter η (Fig.
7 B2). This correlation is not exact, because F_{yw} will depend not only on how much streamflow
8 comes from the upper box, but also on how much of the upper box is young water. That, in
9 turn, will depend on precipitation climatology and the size of the upper box.

10 One can use F_{yw} not only to make comparisons across catchments, but also, in an individual
11 catchment, to compare how the proportions of flow traveling by fast flowpaths change across
12 different flow regimes, as shown in Figs. 8-10 and 12. In turn it may be possible to draw
13 inferences about how catchment processes change with flow regime. In this model, variations
14 in F_{yw} across different flow regimes are strongly correlated with the fractional contributions of
15 the upper box to streamflow (Fig. 13). The slopes and intercepts of the relationships vary
16 among parameter sets, principally reflecting variations in the partitioning parameter η and the
17 sizes of the upper and lower boxes. The strong correlations shown in Fig. 13 are typical.

18 Repeating the analysis shown in Fig. 13 for 200 random model "catchments" (i.e., different
19 random parameter sets) yields an average correlation of over 0.99 (again, with different linear
20 relationships for different parameter values). Of course these results – and, more generally,
21 the interpretation of F_{yw} in terms of upper-box flow – are model-dependent. They are meant
22 to demonstrate only that process inferences can be drawn from F_{yw} , not that these particular
23 inferences should be applied literally to real-world catchments. Indeed one must remember
24 that in the real world there is no "upper box"; it, like all model abstractions, should not be
25 confused with reality.

26 The young water fraction F_{yw} may also be helpful in inferring chemical processes from
27 streamflow concentrations of reactive chemical species. Many reactive species exhibit clear
28 concentration-discharge relationships. Because one can determine how F_{yw} varies, on
29 average, across different ranges of discharge (as demonstrated in Figs. 8-10 and 12), one can
30 potentially construct mixing relationships between F_{yw} and the concentrations of reactive
31 species. If the measurable range of F_{yw} is wide enough, one may even be able to estimate the

1 end-member concentrations corresponding to idealized "young water" ($F_{yw}=1$) and "old
2 water" ($F_{yw}=0$).

3 Figure 14 illustrates a preliminary proof of concept for this approach, based on 20-28 years of
4 weekly precipitation and streamflow samples from three catchments at Plynlimon, Wales
5 (Neal et al., 2011) with contrasting geochemical behavior. I separated the streamflow
6 samples into five discharge ranges (lowest 20 percent, next 20 percent, and so on), then fitted
7 the seasonal chloride concentration cycles in each discharge range and calculated the
8 corresponding young water fractions using the approach outlined in Sect. 3.4 above. I then
9 examined the relationships between these young water fractions and the mean streamwater
10 concentrations of reactive chemical species in each discharge range. Figure 14 shows three
11 different views of how reactive tracer chemistry varies with discharge across the three
12 catchments. The left-hand panels show the average concentrations in each discharge range, as
13 functions of the logarithm of discharge. The middle panels show the same concentrations as
14 functions of the inferred F_{yw} , with the vertical axis at $F_{yw}=0$ indicating the hypothetical old
15 water end-member. The right hand panels show the concentrations plotted against the
16 reciprocal of F_{yw} ; here, the vertical axis at $1/F_{yw}=1$ indicates the hypothetical young water
17 end-member. The gray lines are fitted by hand to indicate general trends, and to suggest
18 potential end-member concentrations.

19 The three catchments are characterized by contrasts in soil hydrology, with the abundance of
20 impermeable gley soils and boulder clay tills increasing in the rank order Hafren < Hore <
21 Tanllwyth. The same rank order is observed in the calculated young water fractions at high
22 flows, reflecting the greater high-flow variability in chloride concentrations at sites with more
23 impermeable soils. The three sites also exhibit contrasting concentration-discharge
24 relationships for nitrate and aluminum (Fig. 14a and d), two solutes that are relatively
25 abundant in near-surface soil solutions. When plotted against the young water fraction,
26 however, these catchment-specific concentration-discharge relationships collapse to single
27 concentration- F_{yw} relationships (Fig. 14b and e) in which the three sites are generally
28 indistinguishable within error. These relationships can be extrapolated to reasonably well-
29 constrained old water end-member concentrations of $\sim 0.1 \text{ mg L}^{-1} \text{ NO}_3\text{-N}$ and $\sim 50 \text{ } \mu\text{g L}^{-1} \text{ Al}$,
30 and to comparably well-constrained young water end-member concentrations of $\sim 0.45 \text{ mg L}^{-1}$
31 $\text{NO}_3\text{-N}$ and $\sim 600 \text{ } \mu\text{g L}^{-1} \text{ Al}$ (Fig. 14c and f). In the case of calcium, the three catchments have
32 markedly different concentration-discharge relationships (Fig. 14g), reflecting differences in

1 the abundance of calcite in their bedrock. As a result, the three catchments have different old
2 water end-member calcium concentrations, ranging from ~ 1 to ~ 4 mg L⁻¹ (Fig. 14h).
3 However, all three streams converge to similar concentrations of ~ 0.5 mg L⁻¹ Ca in the young
4 water end-member (Fig. 14i).

5 It is tempting to interpret the concentration differences between the young and old end-
6 members as reflecting chemical kinetics, but this should be approached with caution. A
7 kinetic interpretation makes sense if the young and old end-members differ only in age (albeit
8 by an unspecified amount since we cannot know how old the "old" end-member is), but not if
9 they differ in other respects as well. At Plynlimon, for example, porewaters in the acidic soil
10 layers have relatively high concentrations of aluminum and transition metals, and relatively
11 low concentrations of base cations and silica, whereas waters infiltrating deep into the
12 fractured bedrock react with calcite and layer lattice silicates and thus become enriched in
13 base cations and silica, and depleted in aluminum and transition metals (Neal et al., 1997).
14 Thus one must also consider the alternative hypothesis that the young end-member represents
15 mostly soil water, and the old end-member represents mostly deeper groundwater, and that
16 the two end-members exhibit different chemistry because of their sources rather than their
17 ages. In this case, the end-member compositions identified through plots like Fig. 14 may
18 help in characterizing the chemistries, and thus localizing the physical sources, of the young
19 and old waters. In this proof-of-concept example, all three catchments appear to have
20 geochemically similar young water end-members, with a composition suggesting a shallow
21 soil source, but each has a different old water end-member, suggesting deeper groundwater
22 sources with differing amounts of carbonate minerals. This is consistent with independent
23 geochemical evidence at Plynlimon (Neal et al., 1997).

24 It is also important to note that if the ideal end-member mixing assumptions hold (i.e., the
25 young and old end-members are invariant, and the mixture undergoes no further chemical
26 reactions), then the mixing relationships in the middle plots of Fig. 14 should be straight lines,
27 and they should extrapolate to physically realistic (non-negative) concentrations at both $F_{yw}=0$
28 and $F_{yw}=1$. To the extent that the mixing relationships are not straight, or imply unrealistic
29 end-members, they indicate that these assumptions are not met.

1 3.7 General observations and caveats

2 It is important to recognize that the inferred young water fractions F_{yw} plotted in Figs. 7-12
3 are not in any way calibrated to the true values determined by age tracking. Nor do they make
4 use of any information about the models that transform precipitation into streamflow (neither
5 their structure, nor their parameter values). Thus there is nothing artifactual about the close
6 correspondence between predicted and observed values of F_{yw} in Figs. 7-12. Instead, these
7 thought experiments provide strong evidence that seasonal tracer cycles can be used to
8 reliably partition streamflow into young and old fractions (F_{yw} and $1-F_{yw}$, respectively), even
9 in catchments that are both nonstationary and spatially heterogeneous, and whose real-world
10 "models" (i.e., whose underlying processes) are poorly understood.

11 When these results are applied in practice, however, one must keep in mind that in contrast to
12 typical field studies, these thought experiments are based on synthetic data sets that are dense
13 (daily measurements for 10 years) and error-free. Furthermore, these thought experiments use
14 a sinusoidal precipitation tracer signal that varies only seasonally, with no confounding
15 variation on shorter or longer time scales. Further benchmark testing will be needed to test
16 the accuracy of F_{yw} estimates derived from shorter, sparser, and messier data sets.

17 One can of course also question the realism of the particular model that I have used for these
18 thought experiments. This model can be calibrated to reproduce the stream discharge with a
19 Nash-Sutcliffe efficiency of better than 0.85 at two of the three sites, but there is no guarantee
20 that it is getting the right answer for the right reasons. All models – whether lumped
21 conceptual models or "physically based" spatially explicit models – necessarily involve
22 approximations and simplifications. In plain language: any model, including this one,
23 incorporates assumptions that are false and are known to be false. One obvious idealization (a
24 less euphemistic word would be *fiction*) is the well-mixed boxes that form the core of most
25 lumped conceptual models, including the model presented here. Assuming that everything in
26 each box is completely mixed or, equivalently, that it is randomly sampled in the outflow –
27 regardless of where it is physically located in the landscape – clearly strains credibility, but
28 this is what typical conceptual models must assume for mathematical convenience. The
29 model presented here is no different.

30 What is different, however, is that here the model is used for purposes that make its literal
31 realism unnecessary. Typical modeling studies draw conclusions about real-world systems
32 from model behavior; thus those conclusions depend critically on the realism of the model.

1 But here, the primary goal is not to test how catchments work, but instead to test specific
2 methods for inferring water ages from complex, nonstationary time series of tracer
3 concentrations. All the model must do is generate outputs with reasonable degrees of
4 complexity and nonstationarity; it is not essential that the model generates these time series by
5 the same mechanisms that real-world catchments do. The only inductive leap is the inference
6 that if a method correctly infers water ages from tracer patterns in these complex,
7 nonstationary time series, it will also correctly infer water ages in complex, nonstationary
8 time series generated by real-world catchments.

9 It is important to highlight an essential difference between the approach developed here and
10 typical studies that infer water ages or transit time distributions from calibrated models (e.g.,
11 Birkel et al., 2011; van der Velde et al., 2012; Heidebüchel et al., 2012; Hrachowitz et al.,
12 2013; Benettin et al., 2015; Benettin et al., 2013). When one draws inferences from a model,
13 their validity depends on whether that model is structurally adequate and whether its
14 parameter values are realistic, both of which are usually in doubt. Here, by contrast, I have
15 developed an inferential method (for estimating the young water fraction F_{yw} from seasonal
16 tracer cycles) that is not drawn from – and thus does not depend on – the model's structure or
17 its parameter values. The model is used only to create synthetic data to test the inferential
18 method.

19 The results reported here, together with those in Paper 1, show that mean transit times
20 (MTT's) cannot be estimated reliably by fitting sine waves to seasonal tracer cycles from
21 nonstationary or spatially heterogeneous catchments. These results do not imply that other
22 methods for estimating MTT's are any better; instead, they imply only that sine wave fitting
23 has been subjected to rigorous benchmark testing and has failed. The other methods have not
24 yet been similarly tested, and it is unclear whether they too will fail. Efforts to fill this
25 knowledge gap are underway. But in the meantime, ignorance is not bliss; one should not
26 simply assume that these other methods work as intended, just because they have not yet been
27 rigorously tested. In that regard, the most general contribution of this analysis is not that it
28 reveals specific problems with MTT estimation from seasonal tracer cycles, or that it
29 demonstrates the reliability of F_{yw} as an alternative metric of catchment transit times, but
30 rather that it illustrates the clarifying power of well-designed benchmark tests.

31

1 **4 Summary and conclusions**

2 The age of streamflow – i.e., the time that has elapsed since it fell as precipitation – is an
3 essential descriptor of catchment functioning with broad implications for runoff generation,
4 contaminant transport, and biogeochemical cycling (Kirchner et al., 2000; McGuire and
5 McDonnell, 2006). The age of streamflow is commonly measured by its mean transit time
6 (MTT), which in turn has often been estimated from the damping of seasonal cycles of
7 chemical and isotopic tracers (such as Cl⁻, δ¹⁸O, or δ²H). In a companion paper ("Paper 1":
8 Kirchner, 2015), I demonstrated that MTT cannot be reliably estimated from seasonal tracer
9 cycles in spatially heterogeneous catchments, and I proposed an alternative water age metric,
10 the young water fraction F_{yw} , which is relatively immune to the errors and biases that afflict
11 the MTT.

12 Here I have explored how catchment nonstationarity affects estimates of MTT and F_{yw} , using
13 simple thought experiments based on a simple two-box conceptual model (Fig. 1), driven by
14 three precipitation time series representing a range of precipitation climatologies (Fig. 2).
15 The model exhibits complex nonstationary behavior (Fig. 3), with striking volatility in tracer
16 concentrations, young water fractions, and mean transit times as the mixing ratio between the
17 upper and lower boxes shifts in response to precipitation events. This mixing ratio is both
18 hysteretic and nonstationary, varying in response both to precipitation forcing and to the
19 antecedent moisture status of the two boxes (Fig. 4).

20 Marginal (time-averaged) age distributions in drainage are skewed toward younger ages than
21 the storage distributions they come from, because storage is flushed more quickly (and thus is
22 younger) during periods of higher discharge (Fig. 5). The age distributions in whole-
23 catchment storage and discharge are approximate power laws, with markedly different slopes
24 (Fig. 5). The age distribution in streamflow becomes increasingly skewed at higher
25 discharges, with a marked increase in the young water fraction and decrease in the mean
26 water age (Fig. 6), reflecting the increased dominance of the upper box at higher flows.
27 Flow-weighted average MTT's are typically close to the steady-state MTT, estimated as the
28 ratio of the total storage to the throughput rate. However, the marginal age distributions are
29 markedly different from the distributions that would be expected in steady state,
30 demonstrating that steady-state approximations are misleading guides to the non-steady-state
31 behavior of the system, *even on average*.

1 Even this simple two-box model exhibits strong equifinality (Fig. B1), with four of its five
2 parameters having virtually no identifiability through hydrograph calibration. However,
3 scaling arguments based on simple perturbation analyses (Sect. 3.2) reveal ratios of
4 parameters that can be constrained through hydrograph calibration (Fig. B2), greatly reducing
5 the equifinality in the parameter space. Unfortunately, water age is primarily controlled by
6 residual storage, which cannot be constrained through hydrograph calibration (Fig. B2).
7 Thus, parameter sets that yield virtually identical hydrographs imply widely differing young
8 water fractions and mean water ages (Fig. B3).

9 The simple two-box model was used to simulate discharge, water ages, and the propagation of
10 seasonal tracer cycles through the catchment, across wide ranges of random parameter sets.
11 MTT's inferred from the damping and phase shift of the seasonal tracer cycles exhibited
12 strong underestimation bias and large scatter (Fig. 7). This result implies that many literature
13 MTT values (and thus also residual storage volumes) may have been underestimated by large
14 factors. By contrast, the seasonal tracer cycles accurately predicted the actual F_{yw} in
15 streamflow, as determined by age tracking within the model (Fig. 7).

16 Flow-weighted fits to the seasonal tracer cycles accurately predicted the flow-weighted
17 average F_{yw} in streamflow, while unweighted fits to the seasonal tracer cycles accurately
18 predicted the unweighted average F_{yw} . The streamflow time series can be separated into
19 distinct flow regimes with their own seasonal tracer cycles (Fig. 8), which accurately reflect
20 the F_{yw} in each flow regime (Figs. 9 and 10). Seasonal tracer cycles also accurately predicted
21 the F_{yw} in the merged streamflow from spatially heterogeneous assemblages of nonstationary
22 model catchments (Fig. 12). Importantly, all of these F_{yw} predictions were really predictions;
23 they were not calibrated in any way.

24 The relationship between F_{yw} and the flow regime reflects how the fluxes from short-term
25 storages vary with hydrologic forcing (Fig. 13). In a preliminary proof of concept (Fig. 14), I
26 showed that one can construct mixing relationships between solute concentrations and F_{yw} 's
27 for discrete flow regimes. From these mixing relationships one can estimate the chemical
28 composition of idealized "young water" and "old water" end-members (Fig. 14).

29 These findings extend the results of Paper 1 by showing that estimates of MTT from seasonal
30 tracer cycles are unreliable under nonstationarity as well as spatial heterogeneity. These
31 findings also extend the results of Paper 1 by showing that F_{yw} can be reliably estimated in
32 nonstationary catchments as well as spatially heterogeneous ones, and can also be reliably

1 estimated for discrete flow regimes. These results further demonstrate that F_{yw} can be reliably
 2 estimated for discrete flow regimes, and can provide helpful insights into the hydrological and
 3 hydrochemical functioning of catchments. Most generally, these results, along with those of
 4 Paper 1, illustrate how well-posed benchmark tests can be essential in clarifying what is
 5 knowable – and, conversely, unknowable – in environmental research.

6

7 **Appendix A: Solution scheme**

8 For simplicity and efficiency, the hydrological model is solved on a fixed daily time step.
 9 This requires some care with the numerics, given the clear (though often overlooked) dangers
 10 in naive forward-stepping simulations of nonlinear equations (Clark and Kavetski, 2010;
 11 Kavetski and Clark, 2010, 2011). Here I use a weighted combination of the trapezoidal
 12 method (which is partly implicit, for enhanced accuracy) and the backward Euler method
 13 (which is fully implicit, for guaranteed stability). The hydrological solution scheme is
 14 illustrated here for the upper box; the lower box is handled analogously. The storage in the
 15 upper box is updated using the following equation:

$$16 \quad S_u(t_{i+1}) - S_u(t_i) = \Delta t \left(P - \rho k_u S_u(t_{i+1})^{b_u} - (1 - \rho) k_u S_u(t_i)^{b_u} \right) \quad , \quad (A1)$$

17 where $S_u(t_i)$ is the storage in the upper box at the beginning of the i^{th} time interval (with length
 18 Δt), $S_u(t_{i+1})$ is the storage at the end of that interval (and thus the beginning of the next), and P
 19 is the average precipitation rate over the interval. Equation (A1) is implicit and nonlinear;
 20 there is no closed-form solution for the future storage $S_u(t_{i+1})$, which instead is found using
 21 Newton's method. The relative dominance of the trapezoidal and backward Euler solutions is
 22 determined by the weighting factor ρ , which takes on values between $\rho=0.5$ (trapezoidal
 23 method) and $\rho=1$ (backward Euler method). The value of ρ in Eq. (A1) is determined for
 24 each time step using the simple stability criterion,

$$25 \quad \rho = \min \left(0.5 + 0.5 \frac{\left(P - k_u S_u(t_i)^{b_u} \right) \Delta t}{\left(P / k_u \right)^{1/b_u} - S_u(t_i)} \quad , \quad 1 \right) \quad , \quad (A2)$$

26 where the numerator represents the amount that S_u would change during one time step if the
 27 instantaneous drainage rate L in Eq. (1) were projected forward in time, and the denominator
 28 represents the difference between S_u 's current value and its equilibrium value at the
 29 precipitation rate P . Equation (A2) says that if the trapezoidal method would move S_u by only

1 a small fraction of the distance to its equilibrium value (at the precipitation rate P), then the
 2 stability advantages of the backward Euler method are unnecessary and the more accurate
 3 trapezoidal method should dominate the solution instead ($\rho \approx 0.5$). On the other hand, if the
 4 trapezoidal method would overshoot the equilibrium value, then $\rho = 1$ and the fully implicit
 5 backward Euler method is used to solve Eq. (A1). The closer the trapezoidal method would
 6 come to overshooting the equilibrium, the larger the value of ρ and the greater the weight that
 7 is given to the backward Euler solution. The guaranteed stability of the backward Euler
 8 method is important when b_u or b_l is large, because the underlying equations can become quite
 9 stiff. After the final value of S_u is determined by Eq. (A1), the drainage from S_u between t_i
 10 and t_{i+1} is determined by mass balance:

$$11 \quad L = P + (S_u(t_i) - S_u(t_{i+1})) / \Delta t \quad , \quad (A3)$$

12 where L is the average drainage rate over the interval Δt between t_i and t_{i+1} .

13 The tracer concentrations are determined under the assumption that each box is well mixed,
 14 implying that individual water parcels within each box do not need to be tracked, and also that
 15 the concentration draining from each box equals the average concentration within the box. I
 16 make the simplifying assumption that each box's inflow and outflow rates (and also inflow
 17 concentrations) are constant over each day. Again taking the upper box as an example, these
 18 assumptions imply that starting from $t = t_i$ the tracer concentration will evolve as

$$19 \quad \frac{dC_u}{dt} = \frac{P(C_P - C_u)}{S_u(t_i) + (P - L)(t - t_i)} \quad , \quad (A4)$$

20 where C_P and C_u are the concentrations in precipitation and the upper box, respectively, and
 21 the denominator expresses how the volume in the box changes with time from its initial value
 22 of $S_u(t_i)$. Integrating Eq. (A4) over an interval Δt yields the concentration updating formula:

$$23 \quad C_u(t_{i+1}) = C_P + (C_u(t_i) - C_P) \left(\frac{S_u(t_i)}{S_u(t_{i+1})} \right)^{(P/(P-L))} \quad , \quad (A5)$$

24 where any quantities that are not shown as functions of time are constant at their average
 25 values over the interval. Equation (A5) could potentially become difficult to compute when P
 26 and L are nearly equal (differing by, say, less than 1 part in 1000), and the power function
 27 approaches its exponential limit. In such cases the change in volume in Eq. (A4) becomes

1 trivially small, and one can replace Eq. (A5) with the more familiar exponential formula for a
 2 well-mixed box of constant volume:

$$3 \quad C_u(t_{i+1}) = C_P + (C_u(t_i) - C_P) \exp(-P \Delta t / S_u) \quad , \quad (A6)$$

4 After the tracer concentrations are updated, the average concentrations in drainage are
 5 calculated by mass balance, as follows:

$$6 \quad C_L = [C_P(t_i) P + C_u(t_i) S_u(t_i) - C_u(t_{i+1}) S_u(t_{i+1})] / L \quad , \quad (A7)$$

7 where C_L is the average concentration in drainage over the time interval between t_i and t_{i+1} .

8 The mean age within each box is modeled analogously to the tracer concentrations, following
 9 the "age mass" concept widely used in groundwater hydrology. Here I will illustrate the
 10 approach using the example of the lower box, since it is the more complex case (for the upper
 11 box, the input age in precipitation is zero, but this is not true for the upper-box drainage that
 12 recharges the lower box). Assuming that the inflow and outflow rates $L(1-\eta)$ and Q_l are
 13 constant over a day, as is the average age $\bar{\tau}_L$ of the inflow from the upper box, the mean age
 14 in the lower box should evolve according to

$$15 \quad \frac{d\bar{\tau}_l}{dt} = \frac{L(1-\eta)(\bar{\tau}_L - \bar{\tau}_l)}{S_l(t_i) + (L(1-\eta) - Q_l)(t - t_i)} + 1 \quad , \quad (A8)$$

16 which is directly analogous to Eq. (A4), except for additional term of +1, which accounts for
 17 the continual aging of the water in the box. The solution to Eq. (A8) is

$$18 \quad \bar{\tau}_l(t_{i+1}) = \bar{\tau}_L + \frac{S_l(t_{i+1})}{2L(1-\eta) - Q_l} + \left(\bar{\tau}_l(t_i) - \bar{\tau}_L - \frac{S_l(t_i)}{2L(1-\eta) - Q_l} \right) \left(\frac{S_l(t_i)}{S_l(t_{i+1})} \right)^{\left(\frac{L(1-\eta)}{L(1-\eta) - Q_l} \right)} \quad , \quad (A9)$$

19 where $\bar{\tau}_l(t_i)$ and $\bar{\tau}_l(t_{i+1})$ are the mean age of the water in the lower box at the beginning and
 20 end of the time interval. Analogously to tracer concentrations, one can calculate the mean age
 21 of the drainage from the box based on the inputs and the change in mean age inside the box,
 22 using conservation of "age mass":

$$23 \quad \bar{\tau}_{Q_l} = [\bar{\tau}_L(t_i) (1-\eta) + \bar{\tau}_l(t_i) S_l(t_i) - (\bar{\tau}_l(t_{i+1}) - \Delta t) S_l(t_{i+1})] / Q_l \quad , \quad (A10)$$

24 where the factor of $-\Delta t$ accounts for the aging of the contents of the box.

25 The approach used here for concentrations and water ages requires the assumption that input
 26 fluxes to each box are constant within each time interval (but constant at their average values,

1 not their initial values). This is a reasonable approximation, particularly when we have no
2 sub-daily precipitation data. And in exchange for this simplifying assumption, equations
3 (A5), (A6), and (A9) provide something important, namely, the exact analytical solution for
4 the evolution of concentration and age during each time interval. Thus these equations
5 directly solve for the correct result even if, for example, an individual day's rainfall is much
6 greater than the total volume of the upper box. The equations above will correctly calculate
7 the consequences of the (potentially many-fold) flushing that occurs in such cases. The
8 approach outlined above also guarantees exact consistency between stocks and fluxes (but
9 note, not in the usual way by updating stocks with fluxes, but rather by calculating output
10 fluxes from inputs and changes in stocks). Readers should keep in mind that all stocks and
11 properties of stocks (i.e., storage volumes, concentrations, and ages) are expressed as the
12 instantaneous values at the beginning of each time interval, and that fluxes and properties of
13 fluxes (i.e., water fluxes and their concentrations and ages) are expressed as averages over
14 each time interval. Otherwise it could be difficult to make sense of the equations above.

15

16 **Appendix B: Equifinality in hydraulic behavior and divergence in travel times**

17 The analysis outlined in Sect. 3.2 implies that approximate equifinality is inevitable, even in
18 such a simple model, because variations in the exponents b_u and b_l and the reference storage
19 levels $S_{u,ref}$ and $S_{l,ref}$ will have nearly offsetting effects on the model's runoff response.
20 Equations (10) and (11) show that, for a given average precipitation forcing, any parameter
21 values for which the partitioning coefficient η and the ratios $S_{u,ref}/b_u$ and $S_{l,ref}/[(1-\eta)b_l]$ are
22 invariant would give nearly equivalent hydrograph predictions, because the hydraulic
23 response timescales of the upper and lower boxes, and their relative contributions to
24 discharge, would be invariant. These conditions can be achieved for widely varying values of
25 the individual parameters b_u , b_l , $S_{u,ref}$, and $S_{l,ref}$.

26 This equifinality problem can be readily visualized by plots like Fig. B1. To generate Fig.
27 B1, I ran the model with Smith River precipitation forcing and the reference parameter set
28 (shown by the red squares in Fig. B1), and used the resulting daily hydrograph (after the spin-
29 up period) as virtual "ground truth" for model calibration. I then ran the model with 1000
30 random parameter sets, and used the Nash-Sutcliffe efficiency (NSE) of the logarithms of
31 discharge to measure how well their hydrographs matched the reference hydrograph (thus the
32 reference hydrograph has NSE=1 by definition). The 50 best-fitting parameter sets, all with

1 NSE \geq 0.98, are shown as dark blue points in Fig. B1. The bottom row of scatterplots shows
2 the conventional "dotty plots". Their flat tops are the hallmark of equifinality, i.e., wide
3 ranges of parameter values give equally good hydrograph predictions (Beven, 2006). Only
4 the partition coefficient η , which performs well across half its range, can be even modestly
5 constrained by calibration. (The other precipitation drivers yield results similar to those
6 shown in Fig. B1.)

7 The other panels of the scatterplot matrix also give important clues to the origins of the
8 observed equifinality. In particular, the best-fitting parameter sets show strong correlations
9 between $S_{u,ref}$ and b_u , and between $S_{l,ref}$ and b_l , as expected from the perturbation analysis
10 presented in Sect. 3.2. Thus good model performance can be obtained across almost the
11 entire range of these parameters, but only for specific parameter combinations. These
12 parameter combinations correspond to "valleys" in the model's response surface, a
13 longstanding problem in model calibration (e.g., Ibbitt and O'Donnell, 1974). The
14 interdependence of the parameters is visually obvious in the scatterplot matrix, but is invisible
15 in the conventional "dotty plots".

16 This information can be exploited to design parameter spaces that are more identifiable
17 through calibration (e.g., Ibbitt and O'Donnell, 1974). An ideal parameter space would be one
18 in which 1) all parameters are highly identifiable, meaning the goodness-of-fit surface is
19 strongly curved along each parameter axis, and 2) in the best-fitting parameter sets, no
20 parameters are strongly correlated with one another. The second of these criteria is necessary
21 (although not sufficient) for the first, as Fig. B1 illustrates. A third criterion is that all
22 parameters that are needed for simulating any quantities of interest must be determined
23 somehow within the parameter space, either individually or through combinations of other
24 parameters. Thus, for example, although the volumes of the boxes ($S_{u,ref}$ and $S_{l,ref}$) are
25 strongly correlated with their exponents (b_u and b_l), the parameter space must allow them to
26 be individually determined, because as Eqs. (12-14) suggest, the mean transit times will be
27 controlled primarily by the volumes alone (not in combination with the exponents), whereas
28 the runoff response will be controlled primarily by the ratios of volumes to exponents (Eqs.
29 10-11). These criteria, plus some trial and error, lead to a more identifiable parameter space,
30 whose five axes are $S_{u,ref}$, $S_{l,ref}$, $S_{u,ref}/(\eta \cdot b_u)$, $S_{l,ref}/b_l$, and η .

31 Figure B2 shows that this parameter space exhibits much less equifinality than the parameter
32 space shown in Fig. B1, although the underlying parameter sets and model simulations are

1 exactly the same. All that has been done is to re-project the parameter space onto a different
2 set of coordinate axes in which the curvature of the goodness-of-fit surface is more clearly
3 visible. Thus, much of the apparent equifinality in the parameter space has been eliminated
4 by simple transformations of variables. These transformations can be designed by eye in this
5 case, because the dimensionality of the original parameter space is low. In higher-dimension
6 parameter spaces, multivariate techniques such as factor analysis may be helpful.
7 Nonetheless, given the obvious utility of this simple correlation analysis and the perturbation
8 analysis of Sect. 3.2, it is surprising that they are not more widely used in hydrological
9 modeling.

10 Despite the improved identifiability of the parameter space, however, it is still not possible to
11 constrain the mean transit time by calibration to the hydrograph. As the bottom row of
12 scatterplots in Figure B2 shows, the mean transit time (MTT) is almost entirely determined by
13 the lower box's reference volume $S_{l,ref}$, as one would expect from Eq. 14. However, as
14 predicted by the perturbation analysis in Sect. 3.2, and as shown by Fig. B2, the runoff
15 response of the model system is essentially independent of $S_{l,ref}$ and therefore cannot be used
16 to constrain it. The runoff response does depend on the ratio of $S_{l,ref}$ to b_l , and thus can be
17 used to constrain that ratio, but it cannot constrain $S_{l,ref}$ by itself, and thus it cannot constrain
18 the MTT. For the young water fraction F_{yw} the outlook is not quite as bleak, because F_{yw} is
19 correlated with the partition coefficient η , which can be constrained somewhat by calibration.
20 As a result, it appears that F_{yw} could potentially be constrained within roughly 1/3 of its full
21 range by parameter calibration to the hydrograph.

22 Figure B3 provides a different visualization of the same equifinality problem. Figure B3
23 shows a two-year excerpt from the simulated time series of streamflows, tracer
24 concentrations, young water fractions, and mean transit times for the reference parameter set
25 (the blue curves), along with the 50 parameter sets that gave the best fit to the reference
26 hydrograph (the gray curves). Because these 50 parameter sets were those that matched the
27 reference hydrograph best, it is unsurprising that the 50 gray hydrographs generally follow the
28 blue reference hydrograph in Fig. B3a. The 50 gray tracer concentration time series also
29 follow the blue reference time series (Fig. B3b), but with somewhat greater variability than
30 the hydrographs, indicating that the parameter values affect the chemographs and the
31 hydrographs in somewhat different ways. But the most striking feature of Fig. B3 is the much
32 greater variability among the young water fractions F_{yw} and (especially) the mean transit

1 times MTT for these same parameter sets (Fig. B3c-d). Although all the parameter sets fit the
2 reference hydrograph nearly perfectly, they vary over a range of 0.3 in F_{yw} (out of a total
3 possible range of 1.0), and over a factor of 9.5 in MTT, on average for the whole time period.
4 Thus these time series demonstrate, consistent with Fig. B2, that there are wide ranges of
5 variability in F_{yw} and especially MTT that cannot be constrained by calibration to the
6 hydrograph.

7

8

9 **Acknowledgements**

10 I thank Scott Jasechko and Jeff McDonnell for the intensive discussions that motivated this
11 analysis, and Markus Weiler and an anonymous reviewer for their comments. I thank the
12 Centre for Ecology and Hydrology for making the Plynlimon data available.

1 **References**

2

3 Benettin, P., van der Velde, Y., van der Zee, S., Rinaldo, A., and Botter, G.: Chloride
4 circulation in a lowland catchment and the formulation of transport by travel time
5 distributions, *Water Resour. Res.*, 49, 4619-4632, doi: 10.1002/wrcr.20309, 2013.

6 Benettin, P., Kirchner, J., Rinaldo, A., and Botter, G.: Modeling chloride transport using
7 travel-time distributions at Plynlimon, Wales, *Water Resour. Res.*, 51, 3259-3276, doi:
8 10.1002/2014WR016600, 2015.

9 Bethke, C. M., and Johnson, T. M.: Groundwater age and groundwater age dating, *Annual*
10 *Review of Earth and Planetary Sciences*, 36, 121-152, doi:
11 10.1146/annurev.earth.36.031207.124210, 2008.

12 Beven, K.: On subsurface stormflow: predictions with simple kinematic theory for saturated
13 and unsaturated flows, *Water Resour. Res.*, 18, 1627-1633, 1982.

14 Beven, K.: A manifesto for the equifinality thesis, *J. Hydrol.*, 320, 18-36, doi:
15 10.1016/j.jhydrol.2005.07.007, 2006.

16 Birkel, C., Soulsby, C., and Tetzlaff, D.: Modelling catchment-scale water storage dynamics:
17 reconciling dynamic storage with tracer-inferred passive storage, *Hydrological Processes*, 25,
18 3924-3936, 2011.

19 Birkel, C., Soulsby, C., Tetzlaff, D., Dunn, S., and Spezia, L.: High-frequency storm event
20 isotope sampling reveals time-variant transit time distributions and influence of diurnal
21 cycles, *Hydrological Processes*, 26, 308-316, doi: 10.1002/hyp.8210, 2012.

22 Botter, G., Bertuzzo, E., and Rinaldo, A.: Transport in the hydrological response: Travel time
23 distributions, soil moisture dynamics, and the old water paradox, *Water Resour. Res.*, 46,
24 W03514, doi: 10.1029/2009WR008371, 2010.

25 Botter, G., Bertuzzo, E., and Rinaldo, A.: Catchment residence and travel time distributions:
26 The master equation, *Geophys. Res. Lett.*, 38, L11403, doi: 10.1029/2011GL047666, 2011.

27 Botter, G.: Catchment mixing processes and travel time distributions, *Water Resour. Res.*, 48,
28 15, W05545, doi: 10.1029/2011wr011160, 2012.

1 Clark, M. P., and Kavetski, D.: Ancient numerical daemons of conceptual hydrological
2 modeling: 1. Fidelity and efficiency of time stepping schemes, *Water Resour. Res.*, 46,
3 W10510, doi: 10.1029/2009wr008894, 2010.

4 Feng, X. H., Faiia, A. M., and Posmentier, E. S.: Seasonality of isotopes in precipitation: A
5 global perspective, *Journal of Geophysical Research-Atmospheres*, 114, D08116, doi:
6 10.1029/2008jd011279, 2009.

7 Harman, C. J.: Time-variable transit time distributions and transport: Theory and application
8 to storage-dependent transport of chloride in a watershed, *Water Resour. Res.*, 51, 1-30, doi:
9 10.1002/2014WR015707, 2015.

10 Heidbüchel, I., Troch, P. A., Lyon, S. W., and Weiler, M.: The master transit time distribution
11 of variable flow systems, *Water Resour. Res.*, 48, W06520, doi: 10.1029/2011WR011293,
12 2012.

13 Hrachowitz, M., Soulsby, C., Tetzlaff, D., Malcolm, I. A., and Schoups, G.: Gamma
14 distribution models for transit time estimation in catchments: Physical interpretation of
15 parameters and implications for time-variant transit time assessment, *Water Resour. Res.*, 46,
16 W10536, doi: 10.1029/2010wr009148, 2010.

17 Hrachowitz, M., Savenije, H., Bogaard, T. A., Tetzlaff, D., and Soulsby, C.: What can flux
18 tracking teach us about water age distribution patterns and their temporal dynamics?, *Hydrol.*
19 *Earth Syst. Sci.*, 17, 533-564, doi: 10.5194/hess-17-533-2013, 2013.

20 Ibbitt, R. P., and O'Donnell, T.: Designing conceptual catchment models for automatic fitting
21 methods, in: *Mathematical Models in Hydrology*, International Association of Hydrological
22 Sciences Publication, 101, Wallingford, U.K., 461-475, 1974.

23 Kavetski, D., and Clark, M. P.: Ancient numerical daemons of conceptual hydrological
24 modeling: 2. Impact of time stepping schemes on model analysis and prediction, *Water*
25 *Resour. Res.*, 46, W10511, doi: 10.1029/2009wr008896, 2010.

26 Kavetski, D., and Clark, M. P.: Numerical troubles in conceptual hydrology: Approximations,
27 absurdities and impact on hypothesis testing, *Hydrological Processes*, 25, 661-670, doi:
28 10.1002/hyp.7899, 2011.

29 Kirchner, J. W., Feng, X., and Neal, C.: Fractal stream chemistry and its implications for
30 contaminant transport in catchments, *Nature*, 403, 524-527, 2000.

1 Kirchner, J. W., Feng, X., and Neal, C.: Catchment-scale advection and dispersion as a
2 mechanism for fractal scaling in stream tracer concentrations, *J. Hydrol.*, 254, 81-100, 2001.

3 Kirchner, J. W.: A double paradox in catchment hydrology and geochemistry, *Hydrological*
4 *Processes*, 17, 871-874, 2003.

5 Kirchner, J. W.: Catchments as simple dynamical systems: catchment characterization,
6 rainfall-runoff modeling, and doing hydrology backward, *Water Resour. Res.*, 45, W02429,
7 doi:10.1029/2008WR006912, 2009.

8 Kirchner, J. W.: Aggregation in environmental systems: Seasonal tracer cycles quantify
9 young water fractions, but not mean transit times, in spatially heterogeneous catchments,
10 *Hydrol. Earth Syst. Sci.*, submitted manuscript, 2015.

11 Kreft, A., and Zuber, A.: On the physical meaning of the dispersion equation and its solutions
12 for different initial and boundary conditions, *Chem. Eng. Sci.*, 33, 1471-1480, 1978.

13 McDonnell, J. J., and Beven, K.: Debates-The future of hydrological sciences: A (common)
14 path forward? A call to action aimed at understanding velocities, celerities and residence time
15 distributions of the headwater hydrograph, *Water Resour. Res.*, 50, 5342-5350, doi:
16 10.1002/2013wr015141, 2014.

17 McGuire, K. J., and McDonnell, J. J.: A review and evaluation of catchment transit time
18 modeling, *J. Hydrol.*, 330, 543-563, 2006.

19 Neal, C., Wilkinson, J., Neal, M., Harrow, M., Wickham, H., Hill, L., and Morfitt, C.: The
20 hydrochemistry of the River Severn, Plynlimon, *Hydrol. Earth Syst. Sci.*, 1, 583-617, 1997.

21 Neal, C., Reynolds, B., Norris, D., Kirchner, J. W., Neal, M., Rowland, P., Wickham, H.,
22 Harman, S., Armstrong, L., Sleep, D., Lawlor, A., Woods, C., Williams, B., Fry, M., Newton,
23 G., and Wright, D.: Three decades of water quality measurements from the Upper Severn
24 experimental catchments at Plynlimon, Wales: an openly accessible data resource for
25 research, modelling, environmental management and education, *Hydrological Processes*, 25,
26 3818-3830, doi: DOI: 10.1002/hyp.8191, 2011.

27 Peters, N. E., Burns, D. A., and Aulenbach, B. T.: Evaluation of high-frequency mean
28 streamwater transit-time estimates using groundwater age and dissolved silica concentrations
29 in a small forested watershed, *Aquatic Geochemistry*, 20, 183-202, 2014.

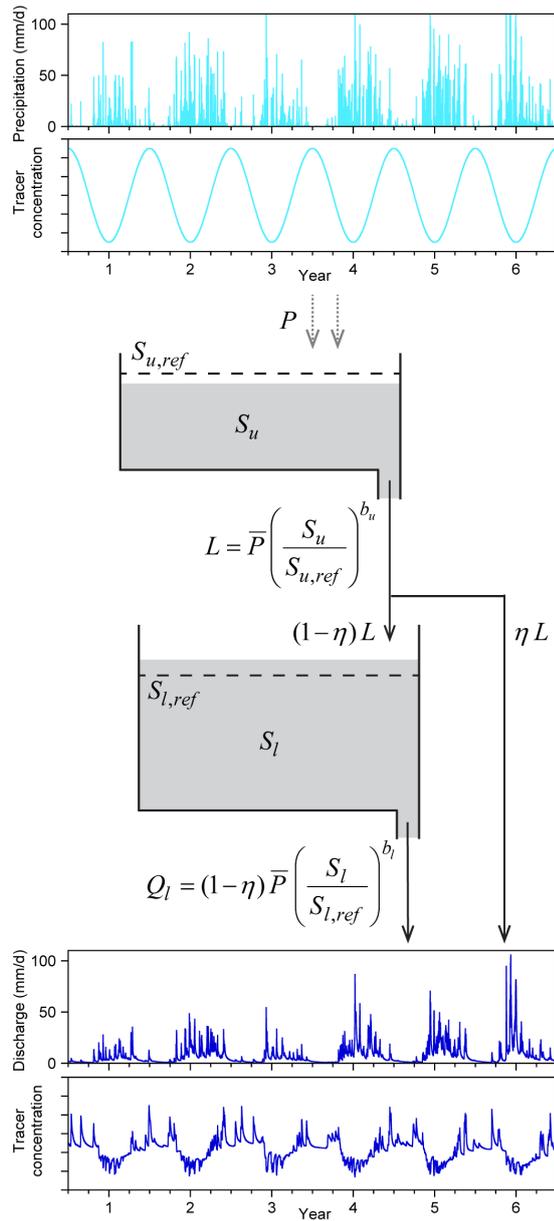
1 Seeger, S., and Weiler, M.: Reevaluation of transit time distributions, mean transit times and
2 their relation to catchment topography, *Hydrol. Earth Syst. Sci.*, 18, 4751-4771, doi:
3 10.5194/hess-18-4751-2014, 2014.

4 Tetzlaff, D., Malcolm, I. A., and Soulsby, C.: Influence of forestry, environmental change and
5 climatic variability on the hydrology, hydrochemistry and residence times of upland
6 catchments, *J. Hydrol.*, 346, 93-111, 2007.

7 Van der Velde, Y., De Rooij, G. H., Rozemeijer, J. C., van Geer, F. C., and Broers, H. P.: The
8 nitrate response of a lowland catchment: on the relation between stream concentration and
9 travel time distribution dynamics, *Water Resour. Res.*, 46, W11534, doi:
10 doi:10.1029/2010WR009105, 2010.

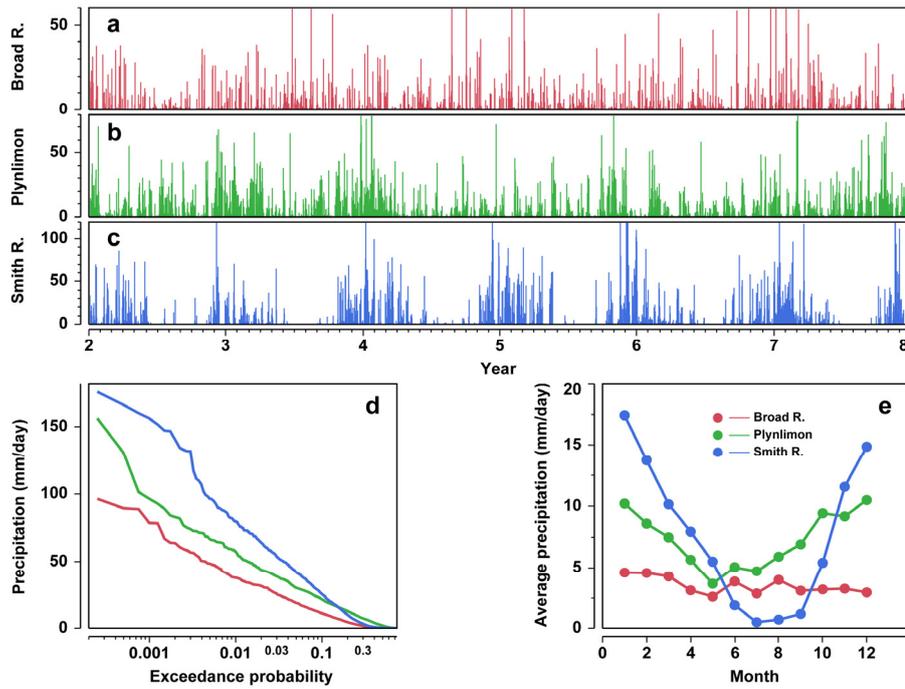
11 van der Velde, Y., Torfs, P. J. J. F., van der Zee, S. E. A. T. M., and Uijlenhoet, R.:
12 Quantifying catchment-scale mixing and its effect on time-varying travel time distributions,
13 *Water Resour. Res.*, 48, W06536, doi: doi:10.1029/2011WR011310, 2012.

14

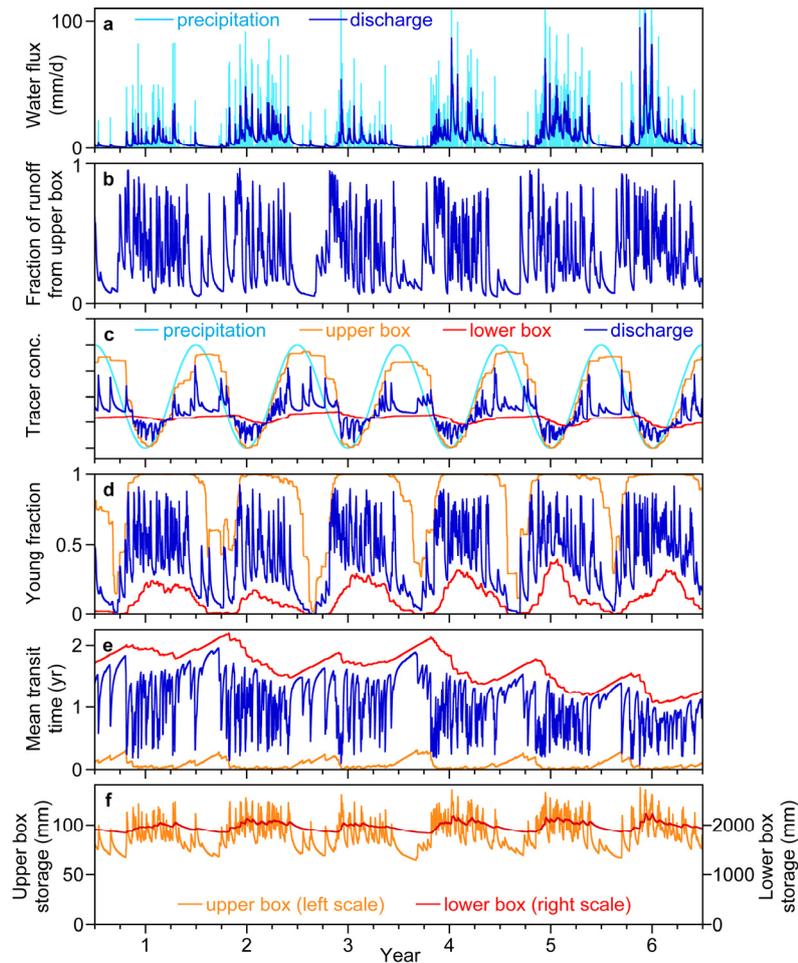


1

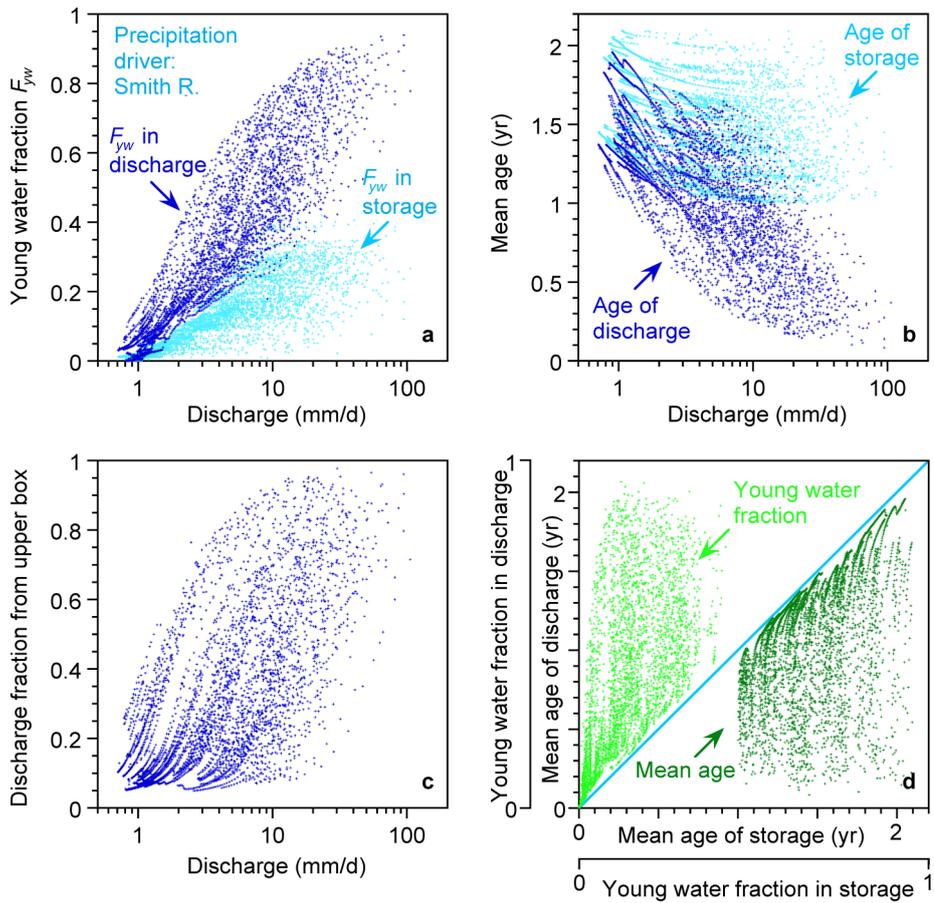
2 Figure 1. Schematic diagram of conceptual model. Drainage from the upper and lower boxes
 3 is determined by power functions of the storage volumes S_u and S_l (depicted by gray shaded
 4 regions) as ratios of the reference storage levels $S_{u,ref}$ and $S_{l,ref}$ (depicted by dashed lines). The
 5 partition coefficient splits the upper box drainage L into direct discharge and infiltration to the
 6 lower box.



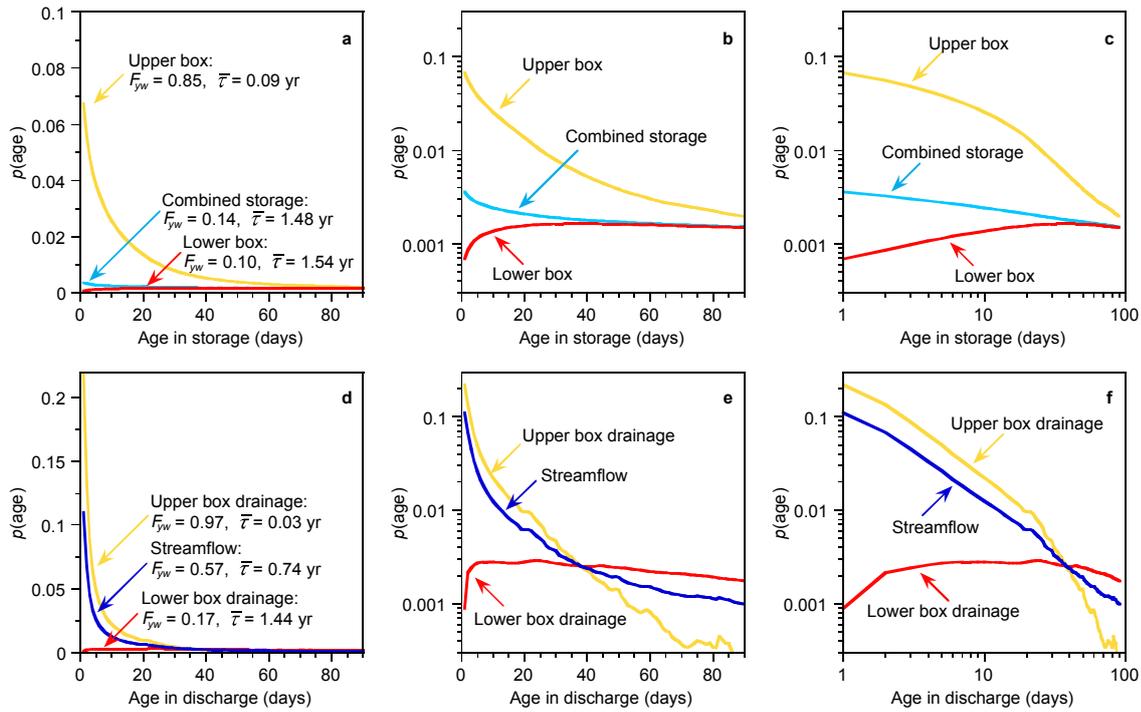
1
 2 Figure 2. Excerpts of daily precipitation records used to drive the model: (a) Broad River,
 3 Georgia, USA (humid temperate climate; Köppen climate zone Cfa) in red, (b) Plynlimon,
 4 Wales (humid maritime climate; Köppen climate zone Cfb) in green, and (c) Smith River,
 5 California, USA (Mediterranean climate; Köppen climate zone Csb) in blue. Axes are
 6 expanded to make typical storms visible; thus the largest storms, some of which extend to
 7 roughly twice the axis limits, are cut off. Exceedance probability plot (d) shows a steeper
 8 magnitude-frequency relationship for Smith River than for the other two records. Monthly
 9 precipitation averages (e) show clear differences in seasonality among the three sites.



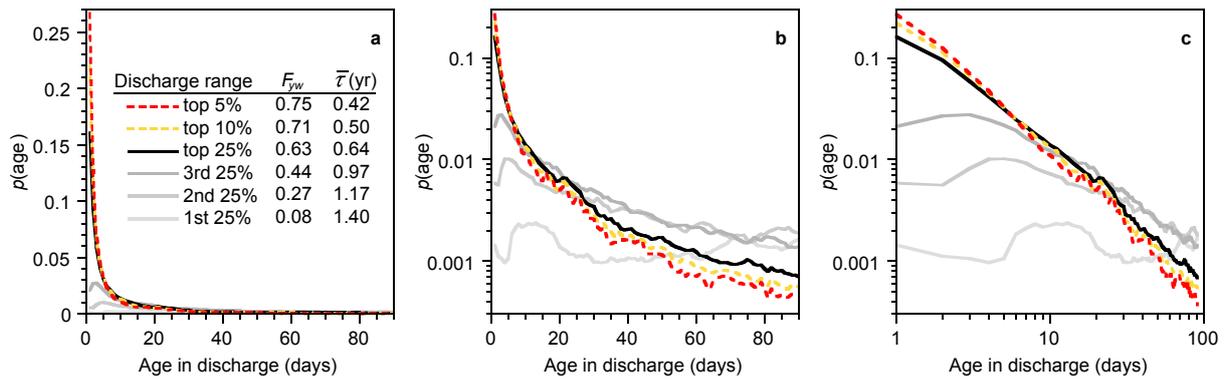
1
 2 Figure 3. Illustrative time series from the two-box model, using the reference parameter set
 3 and the Smith River (Mediterranean climate) precipitation time series. Responses to
 4 precipitation events (a) entail rapid shifts in the proportions of discharge coming from the
 5 upper and lower boxes (b). The smaller, upper box, shown in orange, has a larger young
 6 water fraction (d) and a younger mean age (e) than the larger, lower box, shown in red, and
 7 thus its tracer concentration (c) is less lagged and damped relative to the hypothetical
 8 precipitation concentration, shown by the cosine wave in (c). Mean ages increase (e) and
 9 young water fractions decrease (d), in both boxes, throughout the dry summer periods. The
 10 proportions of streamflow originating from the upper and lower boxes shift dramatically in
 11 response to transient precipitation inputs; thus the tracer concentrations, young water
 12 fractions, and mean ages in discharge (dark blue, panels c-e) vary widely between the time-
 13 varying end members represented by the upper and lower boxes. Storage volumes fluctuate
 14 in a relatively narrow range (f) while discharge varies by orders of magnitude, because the
 15 drainage rates from both boxes are strongly nonlinear functions of storage. Thus both boxes
 16 have sizeable residual storage, which is not drained even under extreme low-flow conditions.



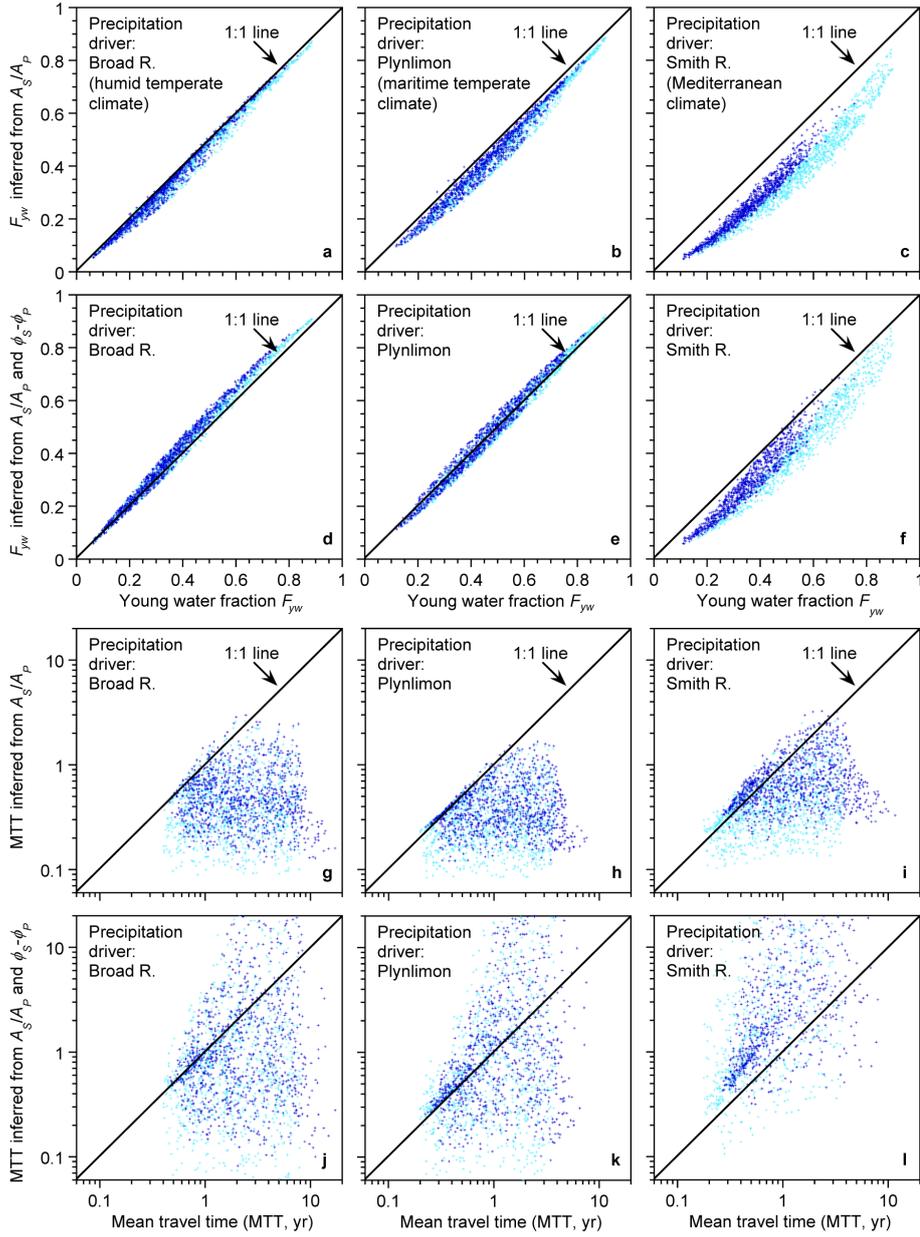
1
 2 Figure 4. Daily values of young water fractions F_{yw} (a) and mean water ages (b) in storage
 3 (light blue) and discharge (dark blue) in the two-box model with reference parameter values
 4 and Smith River (Mediterranean climate) precipitation. The young water fraction and mean
 5 age are both highly scattered functions of discharge (a, b), as is the fractional contribution
 6 from the upper box to streamflow (c), reflecting the effects of variations in antecedent rainfall.
 7 The average age and F_{yw} of water in discharge are strongly biased, and highly scattered,
 8 measures of the same quantities in storage (d).



1
2 Figure 5. Marginal (time-averaged) age distributions in storage (a-c) and drainage (d-f) in the
3 reference case simulation (Fig.3), shown on linear (a, d), log-linear (b, e), and double-log (c,
4 f) axes. Distributions in drainage (lower plots) are skewed toward younger ages than the
5 storage distributions that they come from (upper plots). This arises, even though drainage is
6 not age-selective, because storage is flushed more quickly (and thus is younger) during
7 periods of higher discharge. Age distributions in the upper box, combined storage, and
8 streamflow are more skewed than exponentials (i.e., they are upward-curving in the middle
9 plots). The age distributions in the combined storage and streamflow (blue lines) are
10 approximate power laws; i.e., they are nearly straight in the right-hand plots, with markedly
11 different power-law slopes. The light blue line in the upper plots shows the age distribution
12 of the combined upper and lower boxes, which resembles the age distribution of the lower
13 box because the reference parameter values imply that the lower box comprises about 95
14 percent of total storage. However, direct drainage from the upper box comprises 50 percent
15 of streamflow; thus the streamflow age distribution (shown by the dark blue line in lower
16 plots) reflects the strong skew of the upper box age distribution. Although both boxes are
17 well mixed and have nearly constant volumes, the age distribution of discharge clearly differs
18 from the distribution that would be expected in steady state, which would be exponential in
19 short time.

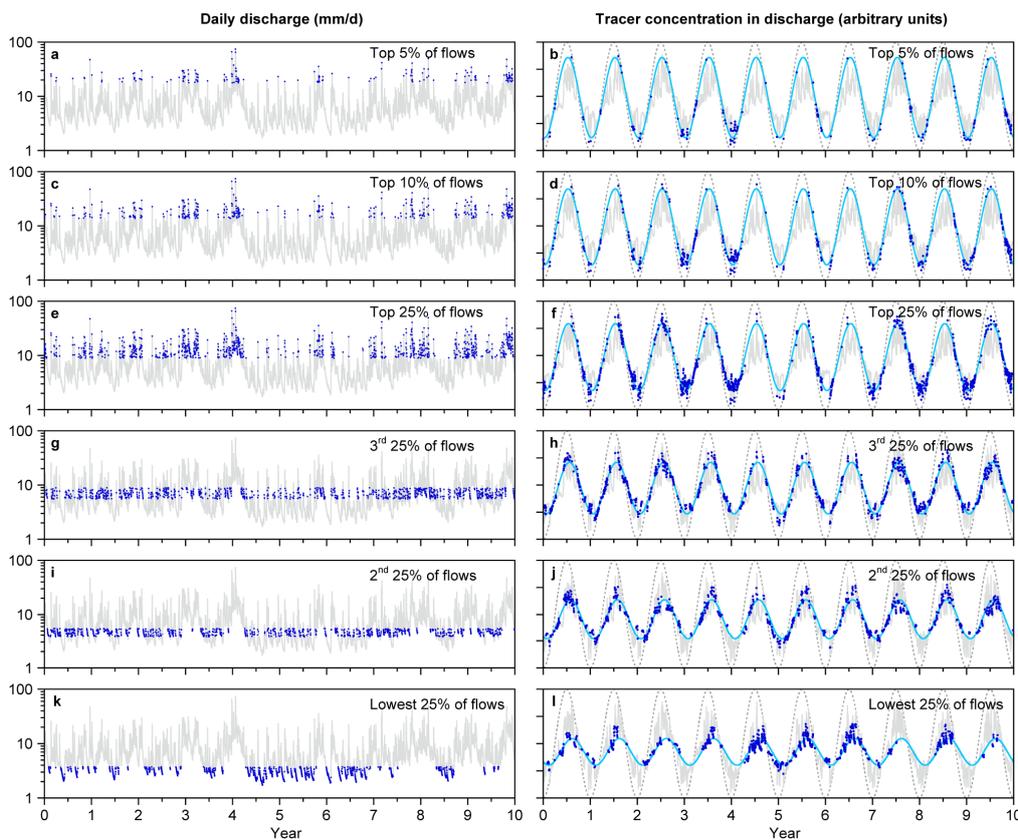


1
 2 Figure 6: Marginal (time-averaged) transit-time distributions (TTD's) for selected ranges of
 3 daily discharges in the two-box model, with the reference parameter set and Smith River
 4 (Mediterranean climate) precipitation forcing, on linear (a), log-linear (b), and double-log (c)
 5 axes. The TTD becomes increasingly skewed at higher discharges (a), with a marked increase
 6 in the young water fraction F_{yw} and decrease in the mean water age $\bar{\tau}$. For the upper half of
 7 all discharges, the age distribution is upward-curving on log-linear axes (b), implying that it is
 8 more skewed than exponential. Discharges in the top 25% and above have approximately
 9 power-law age distributions, plotting as nearly straight lines on double-log axes (c).



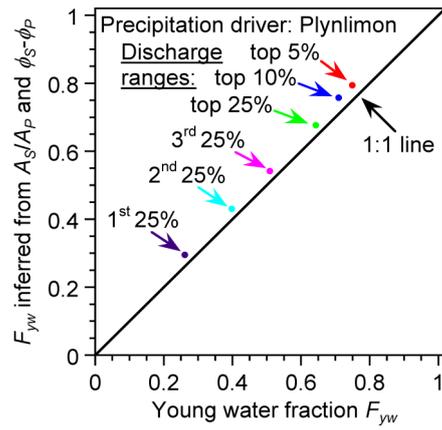
1
2 Figure 7. Young water fractions (F_{yw} , left panels) and mean transit times (MTT, right panels
3 – note log scale) in streamflow from the two-box model. Upper panels compare the average
4 F_{yw} in discharge, determined by age tracking within the model (on the horizontal axes) with
5 the seasonal tracer cycle amplitude ratio A_S/A_P (panels a-c), and with F_{yw} inferred from the
6 tracer cycle amplitude ratio A_S/A_P and phase shift $\phi_S-\phi_P$ (panels d-f). Lower panels compare
7 the average MTT in discharge (again from age tracking) with MTT inferred from the tracer
8 amplitude ratio (panels g-i) and from amplitude ratio and phase shift (panels j-l). Light blue
9 points show flow-weighted average F_{yw} 's and MTT's for each simulation, compared to
10 estimates from flow-weighted fits to seasonal tracer cycles. Dark blue points show un-
11 weighted average F_{yw} 's and MTT's, compared to estimates from un-weighted fits to seasonal

1 tracer cycles. Panels show results from 1000 random parameter sets and three contrasting
 2 precipitation drivers: Broad River (humid, temperate, with very little seasonality), Plynlimon
 3 (wet maritime climate with slight seasonality), and Smith River (Mediterranean climate with
 4 pronounced winter-wet, summer-dry seasonality). Seasonal tracer cycle amplitudes generally
 5 predict the average young water fraction, although they exhibit some systematic bias under
 6 strongly seasonal precipitation regimes like Smith River, where seasonal cycles in
 7 precipitation volume are correlated with seasonal cycles in tracer concentration. By contrast,
 8 mean transit time estimates from seasonal tracer cycles are highly unreliable in all
 9 precipitation regimes.

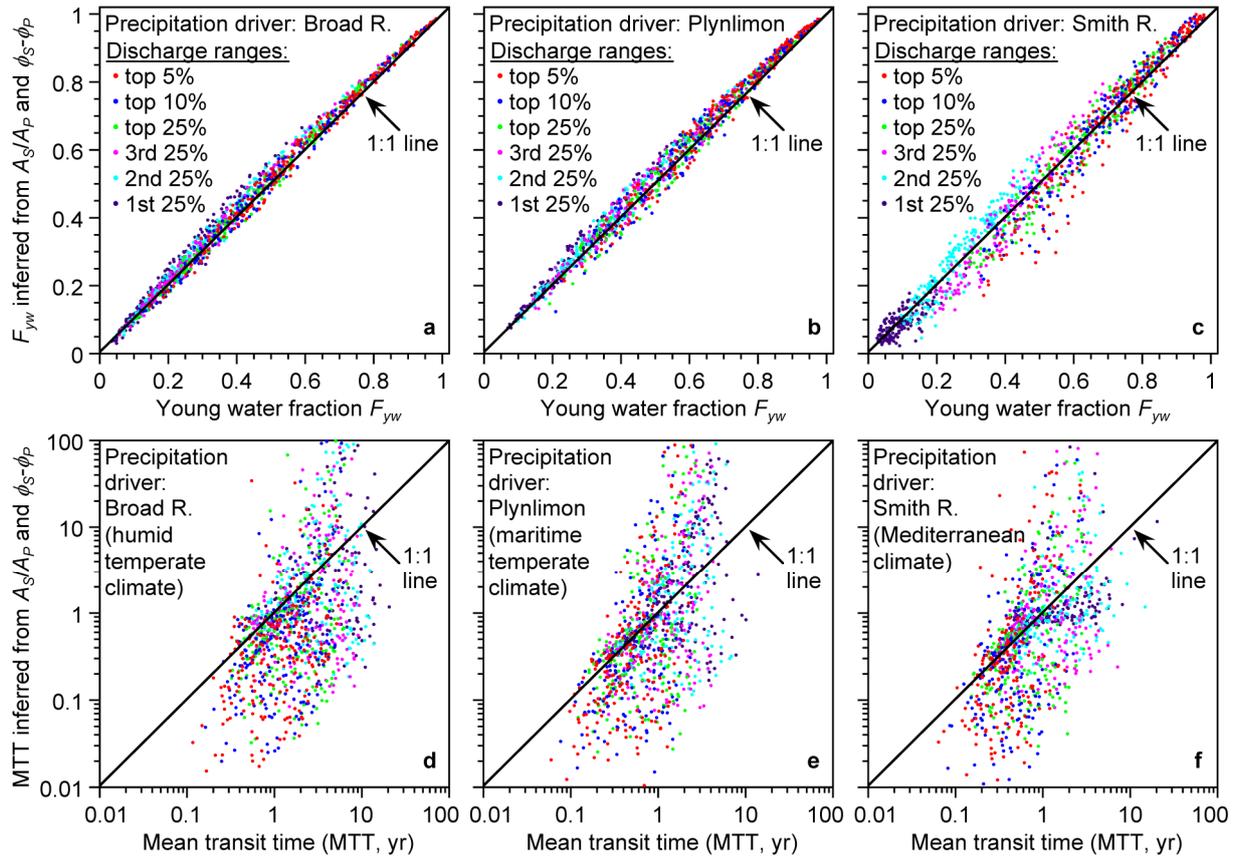


10

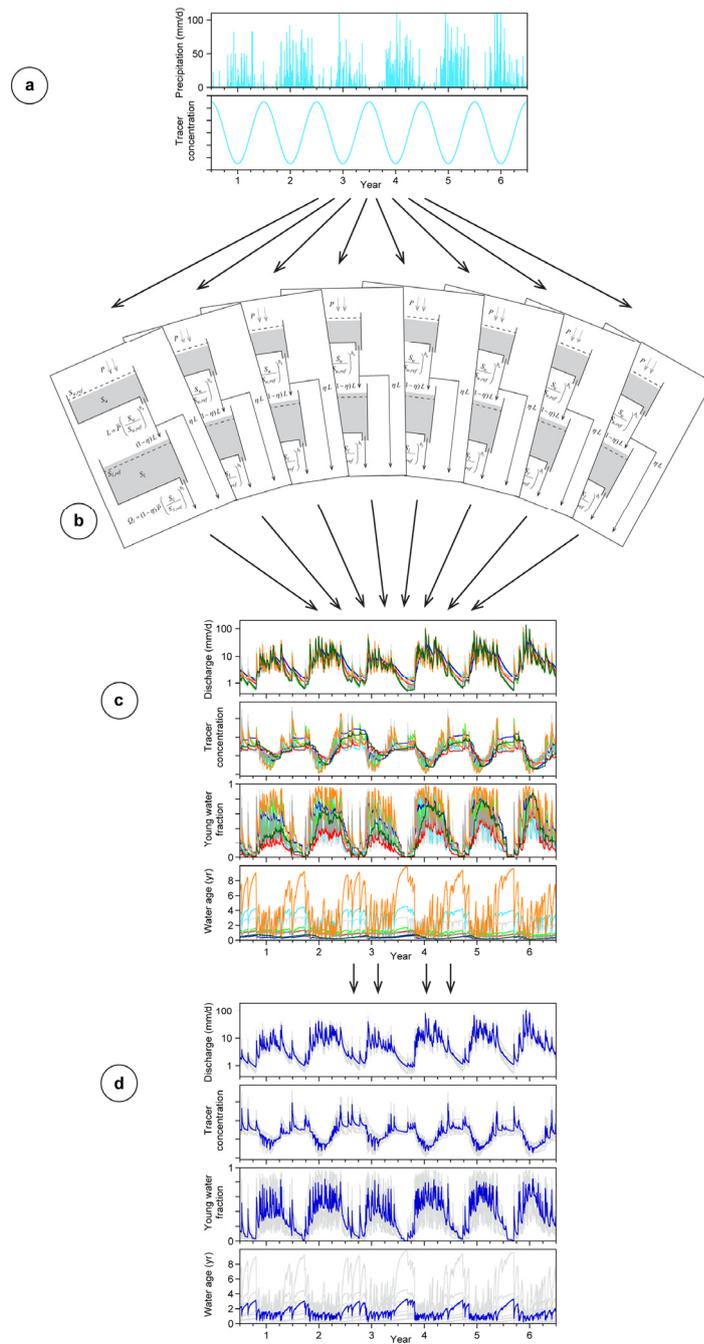
11 Figure 8. Daily discharges (left panels) and tracer concentrations (right panels) in streamflow
 12 from two-box model with reference parameter values and Plynlimon precipitation forcing.
 13 Individual discharge ranges and corresponding tracer concentrations are highlighted in dark
 14 blue. In right-hand panels, precipitation tracer concentrations are shown by dashed gray lines
 15 and best-fit sinusoidal fits to streamflow tracer concentrations are shown in light blue. At
 16 higher discharges, tracer cycles are less damped and less phase-shifted, indicating greater
 17 fractions of young water in streamflow.



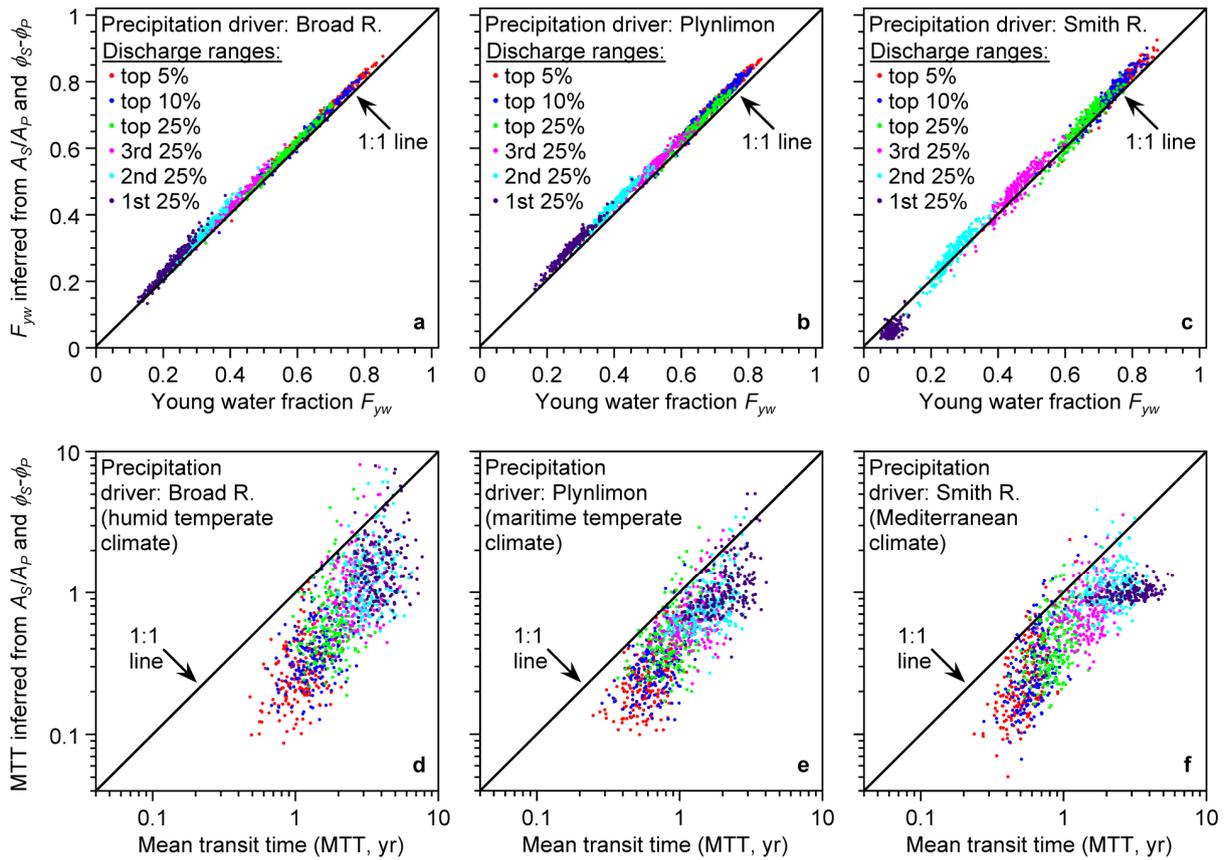
- 1
- 2 Figure 9. Time-averaged, flow-specific young water fractions F_{yw} for the six discharge ranges
- 3 shown in Fig. 8, measured by age tracking in the model (with Plynlimon precipitation forcing
- 4 and the reference parameter set), compared to F_{yw} values estimated from the amplitude ratios
- 5 A_S/A_P and phase shifts $\phi_S - \phi_P$ of the tracer cycles shown in Fig. 8.



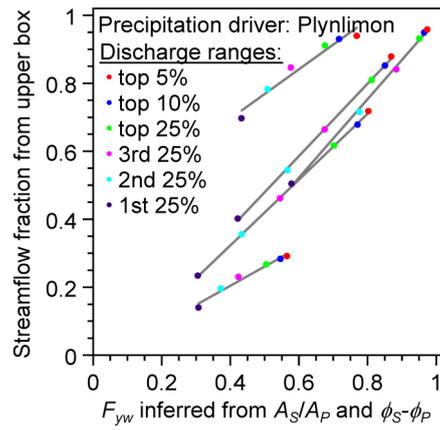
1
2 Figure 10. Young water fractions (F_{yw}) and mean transit times (MTT) in separate discharge
3 ranges in streamflow from two-box model. Upper panels compare the time-averaged, flow-
4 specific F_{yw} for each discharge range (measured by age tracking in the model) with F_{yw} values
5 estimated from the amplitude ratios A_S/A_P and phase shifts $\phi_S-\phi_P$ of the best-fit tracer cycle
6 sinusoids in those discharge ranges (analogously to Fig. 8) using Eqs. (10)-(11) and (13)-(14)
7 of Paper 1. Similar results (not shown) are also obtained for flow-weighted F_{yw} and flow-
8 weighted tracer cycle sinusoids. Results obtained from tracer cycle amplitude alone (without
9 phase information) are also similar, except in some cases where the amplitude ratio is small
10 (particularly with Smith River precipitation forcing). Lower panels compare the MTT,
11 determined by age tracking, with the MTT inferred from tracer amplitude ratios and phase
12 shifts using Eqs. (10)-(11) from Paper 1. Each panel shows results from 200 random
13 parameter sets and three contrasting precipitation drivers: Broad River (humid, temperate,
14 with very little seasonality), Plynilimon (wet maritime climate with slight seasonality), and
15 Smith River (Mediterranean climate with pronounced winter-wet, summer-dry seasonality).
16 Tracer cycle amplitudes and phases generally predict the young water fractions in each
17 discharge range, although with some modest scatter. Mean transit time estimates, by contrast,
18 are highly unreliable, exhibiting large scatter (note log scales).



1
 2 Figure 11. Scheme for simulating spatially heterogeneous catchments with nonstationary
 3 tributary subcatchments. A single precipitation time series (a) is used to drive eight copies of
 4 the model representing eight tributary streams (b), each with a different set of random
 5 parameter values. Streamflows, tracer concentrations, young water fractions, and water ages
 6 from these eight nonstationary tributaries (c, with each color representing a separate tributary
 7 stream) are mass-averaged to determine the time series that would be observed in the merged
 8 streamflow (d, with blue lines showing the merged streamflow and gray lines showing the
 9 tributaries).

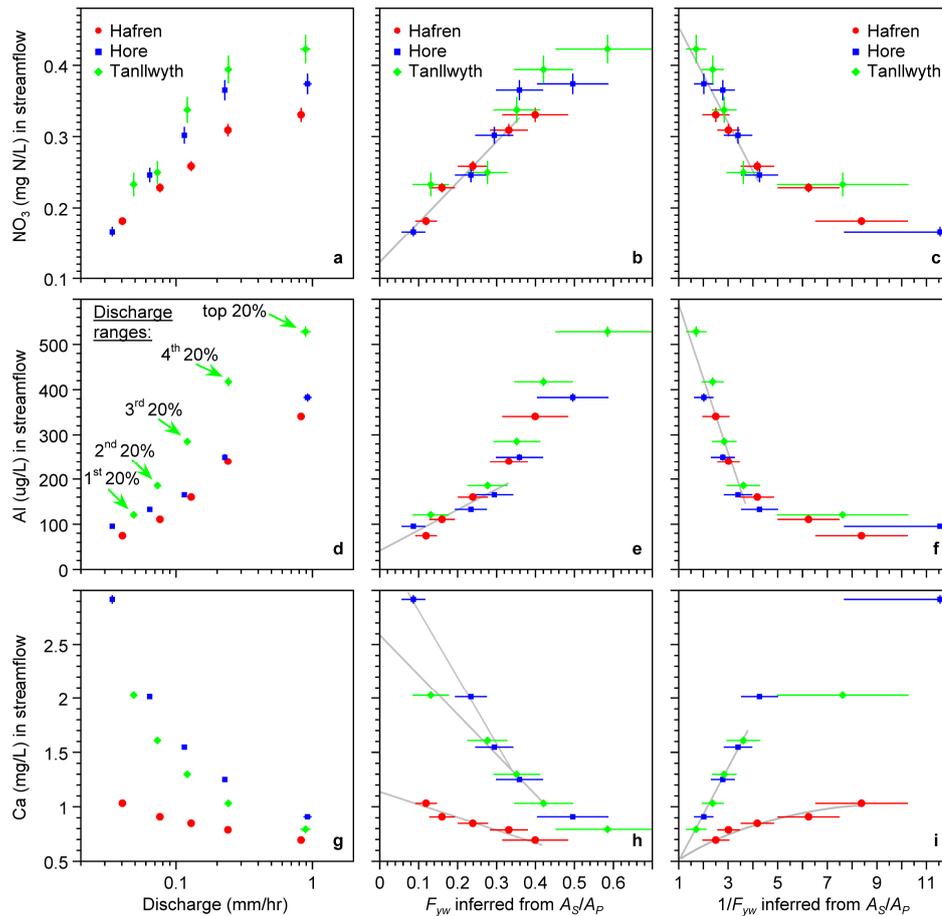


1
2 Figure 12. Actual and inferred young water fractions (F_{yw} , top panels) and mean transit times
3 (MTT, bottom panels) in separate discharge ranges, under combined effects of nonstationarity
4 and spatial heterogeneity. Panels show results for 200 synthetic catchments, each consisting
5 of 8 copies of the two-box model with independent random parameter sets (Fig. 11). Upper
6 panels compare average F_{yw} 's with F_{yw} 's predicted from amplitudes and phases of best-fit
7 tracer cycle sinusoids for each discharge range (e.g., Fig. 8) using Eqs. (10)-(11) and (13)-
8 (14) of Paper 1. Similar results (not shown) are also obtained for flow-weighted F_{yw} 's and
9 flow-weighted tracer cycle sinusoids. Results obtained from tracer cycle amplitude alone
10 (without phase information) are also similar, but exhibit slightly greater bias. Lower panels
11 compare MTT with MTT predicted from tracer amplitude ratios and phase shifts using Eqs.
12 (10)-(11) from Paper 1. Seasonal tracer cycle amplitudes and phases accurately predict young
13 water fractions in separate flow regimes; the corresponding estimates of mean transit times
14 exhibit substantial bias and scatter.

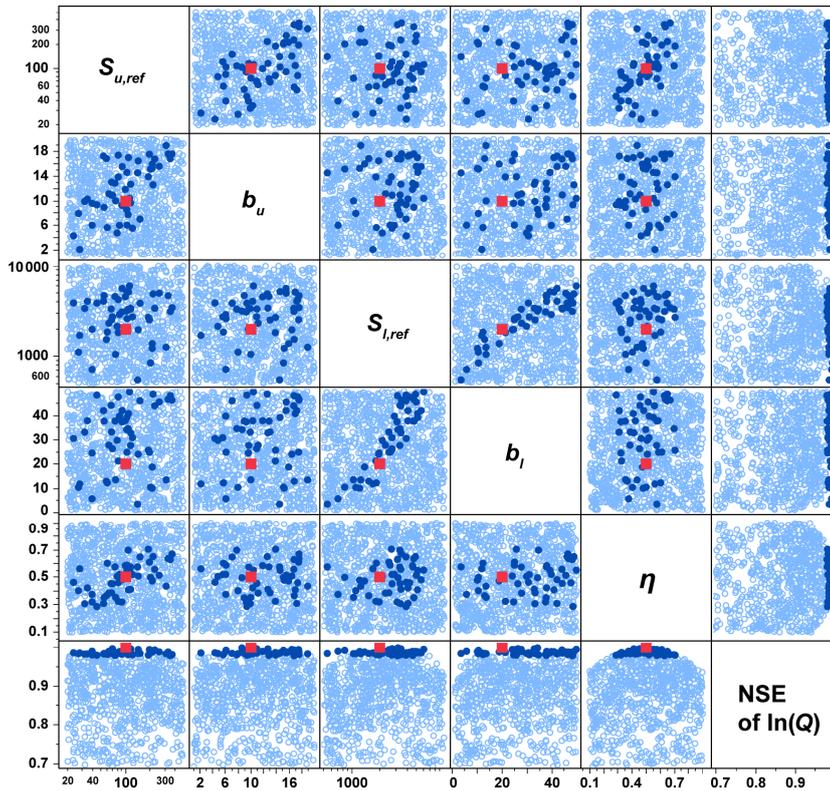


1

2 Figure 13. Correlations between flow-weighted young water fractions F_{yw} and fractional
 3 contributions of the upper box to streamflow across different discharge ranges, for five
 4 parameter sets illustrating the diversity of relationships that can arise in the model. The upper
 5 box contribution is strongly correlated with F_{yw} in all cases, although the slopes and the
 6 intercepts vary among parameter sets.

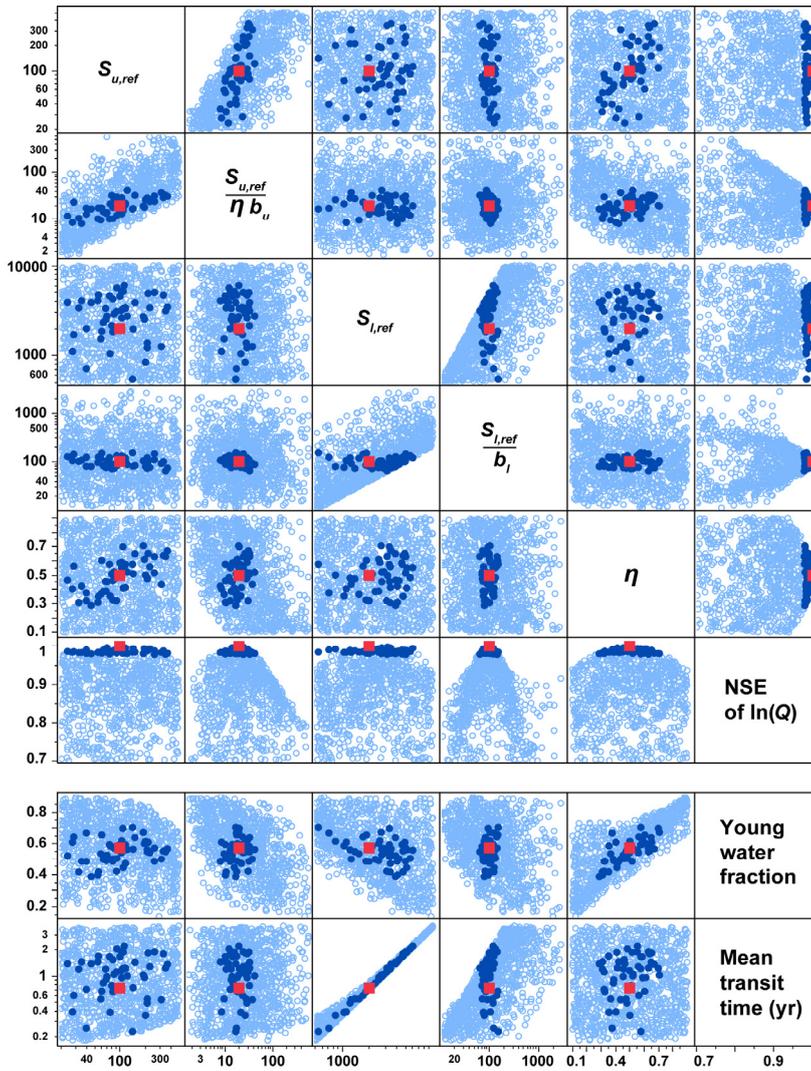


1
2 Figure 14. Concentrations of reactive chemical species as functions of discharge (left panels),
3 young water fractions (middle panels), and reciprocal young water fractions (right panels) for
4 streams draining three contrasting catchments at Plynlimon, Wales. Symbols show means for
5 20% intervals of each catchment's discharge distribution, and error bars indicate ±1 standard
6 error. Gray lines are drawn by hand to indicate general trends. Concentration-discharge
7 relationships in nitrate and aluminum differ among the three catchments (a and d), but
8 collapse to single concentration- F_{yw} relationships (b-c and e-f). These concentration- F_{yw}
9 relationships extrapolate to broadly consistent old water end-members ($F_{yw}=0$, panels b and e)
10 and young water end-members ($F_{yw}=1$, panels d and f). Calcium follows different
11 concentration- F_{yw} relationships in the three streams, which extrapolate to three different old
12 water end-members (h) but roughly the same young water end-member (i).



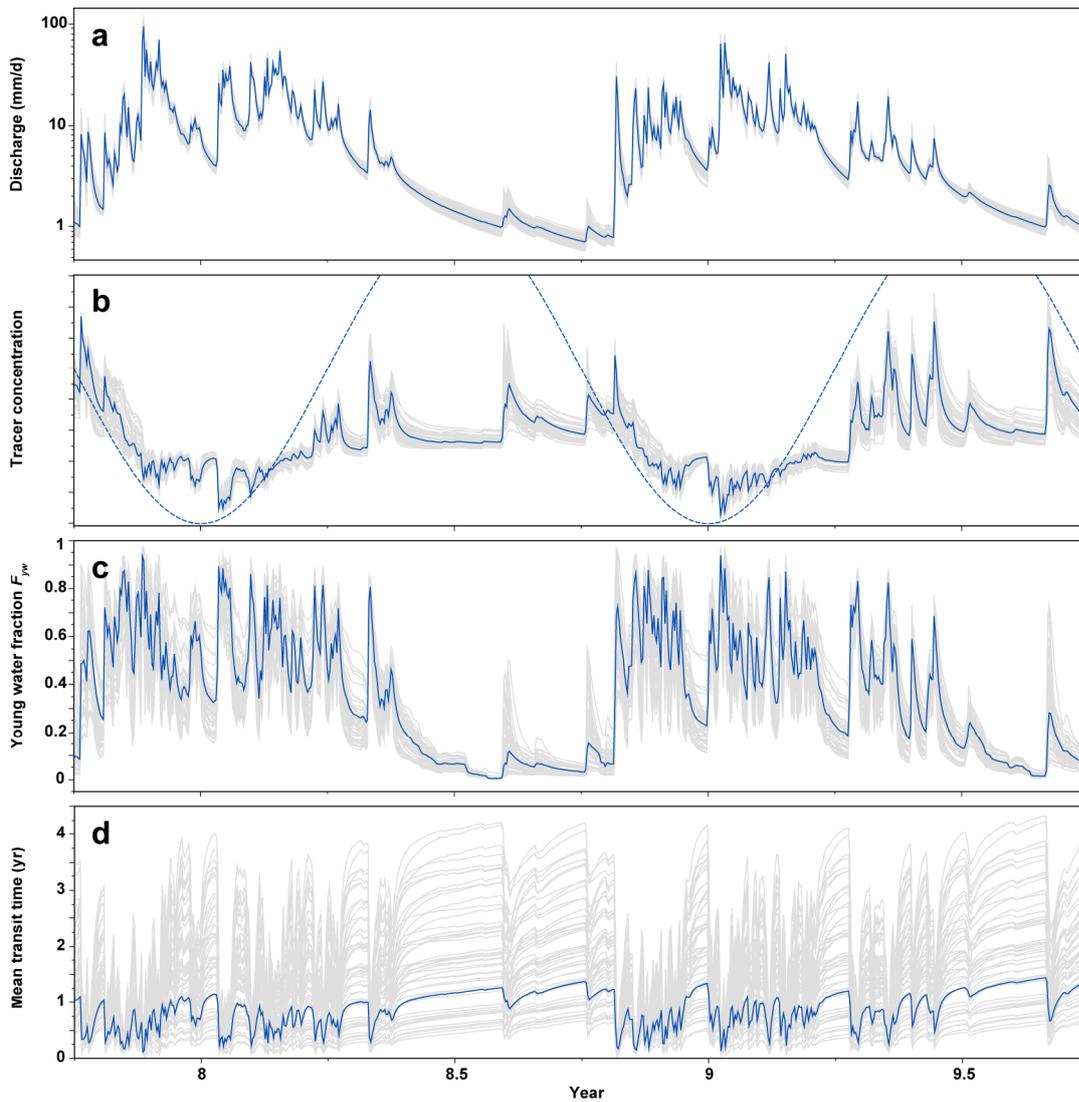
1

2 Figure B1. Equifinality in discharge predictions. The scatterplot matrix shows relationships
 3 among 1000 random parameter sets and the Nash-Sutcliffe efficiency (NSE) of discharge time
 4 series driven by Smith River (Mediterranean climate) precipitation forcing. The red square
 5 indicates the "reference" parameter set that was used to generate the discharge time series that
 6 the other parameter sets were tested against; these reference parameters thus correspond to
 7 NSE=1.00 by definition. The dark blue dots show the best-fitting 50 (or 5%) of the parameter
 8 sets, all with $NSE \geq 0.98$. Excellent discharge predictions can be obtained across almost the
 9 full range of all five model parameters, except the partition coefficient η , which performs well
 10 across only about half its range. The dark blue dots show clear correlations between the
 11 reference storage levels in each box ($S_{u,ref}$, $S_{l,ref}$) and the corresponding drainage function
 12 exponents (b_u , b_l); these correlations delimit regions with nearly constant hydraulic response
 13 time scales, as defined by Eqs. (10)-(11).



1

2 Figure B2. Equifinality partly cured by parameter transformations. The scatterplot matrix
 3 shows relationships among 1000 random parameter sets and the Nash-Sutcliffe efficiency
 4 (NSE) of discharge time series driven by Smith River (Mediterranean climate) precipitation
 5 forcing, along with two key model outputs, the young water fraction and mean transit time in
 6 discharge (bottom two rows). As in Fig. B1, the red square indicates the "reference"
 7 parameter set that was used to generate the discharge time series that the other parameter sets
 8 were tested against; these reference parameters thus correspond to NSE=1.00 by definition.
 9 The dark blue dots show the best-fitting 50 (or 5%) of the parameter sets, all with NSE \geq 0.98.
 10 In contrast to Fig. B1, three of the five parameters can be constrained by calibration against
 11 discharge (as shown by the clear peaks in NSE), and none of the parameters are strongly
 12 correlated with one another. However, the two reference storage volumes $S_{u,ref}$ and $S_{l,ref}$
 13 remain poorly constrained. The mean transit time is determined almost entirely by $S_{l,ref}$, so it
 14 cannot be constrained by parameter calibration against the streamflow hydrograph.



1
 2 Figure B3. Excerpts from time series of discharge, tracer concentrations, young water
 3 fractions, and mean travel times in the two-box model with Smith River (Mediterranean
 4 climate) precipitation forcing and the reference parameter set (the dark lines, for the
 5 parameter values shown by the red squares in Figs. B1 and B2) and the 50 parameter sets that
 6 come closest to matching the reference discharge time series (the light gray lines, for the
 7 parameter sets shown by the solid blue dots in Figs. B1 and B2). The 50 gray hydrographs
 8 (panel a) cluster closely around the blue hydrograph (which is unsurprising because they have
 9 been selected to do so). The 50 gray tracer concentration curves (panel b) also generally
 10 follow the blue curve (the precipitation tracer sinusoid is shown for comparison by the dashed
 11 line). By contrast, the young water fraction F_{yw} (panel c) and mean transit time (panel d) are
 12 much more variable; the gray curves vary by an average range of 0.3 in F_{yw} and a factor of 9.5
 13 in mean transit time.