The Plumbing of Land Surface Models: Is Poor Performance a Result of Methodology or Data Quality?

NED HAUGHTON,^a GAB ABRAMOWITZ,^a ANDY J. PITMAN,^a DANI OR,^b MARTIN J. BEST,^c HELEN R. JOHNSON,^c GIANPAOLO BALSAMO,^d AARON BOONE,^e MATTHIAS CUNTZ,^f BERTRAND DECHARME,^e PAUL A. DIRMEYER,^g JAIRUI DONG,^h MICHAEL EK,^h ZICHANG GUO,^g VANESSA HAVERD,ⁱ BART J. J. VAN DEN HURK,^j GREY S. NEARING,^k BERNARD PAK,¹ JOE A. SANTANELLO JR.,^k LAUREN E. STEVENS,¹ AND NICOLAS VUICHARD^m

^a ARC Centre of Excellence for Climate Systems Science, Sydney, New South Wales, Australia ^b Department of Environmental Systems Science, ETH Zurich, Zurich, Switzerland ^c Met Office, Exeter, United Kingdom ^d ECMWF, Reading, United Kingdom ^c CNRM-GAME, Météo-France, Toulouse, France

^f Helmholtz Centre for Environmental Research (UFZ), Leipzig, Germany

^g Center for Ocean–Land–Atmosphere Studies, George Mason University, Fairfax, Virginia

^hNOAA/NCEP/EMC, College Park, Maryland

ⁱ Oceans and Atmosphere, CSIRO, Canberra, Australian Capital Territory, Australia

^j Royal Netherlands Meteorological Institute (KNMI), De Bilt, Netherlands

Hydrological Sciences Laboratory, NASA GSFC, Greenbelt, Maryland

¹Oceans and Atmosphere, CSIRO, Aspendale, Victoria, Australia ^m Laboratoire des Sciences du Climat et de l'Environnement, UMR 8212, IPSL-LSCE, CEA-CNRS-UVSQ, Gif-sur-Yvette, France

(Manuscript received 16 September 2015, in final form 27 January 2016)

ABSTRACT

The Protocol for the Analysis of Land Surface Models (PALS) Land Surface Model Benchmarking Evaluation Project (PLUMBER) illustrated the value of prescribing a priori performance targets in model intercomparisons. It showed that the performance of turbulent energy flux predictions from different land surface models, at a broad range of flux tower sites using common evaluation metrics, was on average worse than relatively simple empirical models. For sensible heat fluxes, all land surface models were outperformed by a linear regression against downward shortwave radiation. For latent heat flux, all land surface models were outperformed by a regression against downward shortwave radiation, surface air temperature, and relative humidity. These results are explored here in greater detail and possible causes are investigated. It is examined whether particular metrics or sites unduly influence the collated results, whether results change according to time-scale aggregation, and whether a lack of energy conservation in flux tower data gives the empirical models an unfair advantage in the intercomparison. It is demonstrated that energy conservation in the observational data is not responsible for these results. It is also shown that the partitioning between sensible and latent heat fluxes in LSMs, rather than the calculation of available energy, is the cause of the original findings. Finally, evidence is presented that suggests that the nature of this partitioning problem is likely shared among all contributing LSMs. While a single candidate explanation for why land surface models perform poorly relative to empirical benchmarks in PLUMBER could not be found, multiple possible explanations are excluded and guidance is provided on where future research should focus.

DOI: 10.1175/JHM-D-15-0171.1

Corresponding author address: Ned Haughton, ARC Centre of Excellence for Climate Systems Science, University of New South Wales, Level 4, Matthews Building, Sydney NSW 2052, Australia. E-mail: ned@nedhaughton.com



FIG. 1. The locations of the 20 flux tower sites in the PLUMBER experiment. The IGBP vegetation type is represented by color and the numbers indicate the years of data used in the PLUMBER experiment. Site data are given in Table 1.

1. Introduction

The assessment and intercomparison of land surface models (LSMs) has evolved from simple, site-based synthetic experiments in the absence of constraining observational data (Henderson-Sellers et al. 1996; Pitman et al. 1999) to targeted comparisons of process representation (e.g., Koster et al. 2006; Guo et al. 2006) and global-scale experiments (Dirmeyer et al. 1999; Koster et al. 2004; Seneviratne et al. 2013). This history is detailed in Pitman (2003), van den Hurk et al. (2011), Dirmeyer (2011), and Best et al. (2015). Recently, Best et al. (2015) noted that throughout this history, model performance has been assessed by direct comparison with observational products or other LSMs. They argued that without a mechanism to define appropriate levels of performance in a given metric, simple comparisons of this nature are not sufficient to gauge whether models are performing well or not.

The Protocol for the Analysis of Land Surface Models (PALS) Land Surface Model Benchmarking Evaluation Project (PLUMBER; Best et al. 2015) was constructed to undertake a multimodel examination of LSMs and to focus on defining benchmarks for model performance, rather than simply comparing LSMs and observations. PLUMBER examined the performance of 13 LSMs consisting of variants from eight distinct models at 20 flux tower sites covering a wide variety of biomes (see Fig. 1 and Tables 1, 2). Part of the assessment of performance used four common metrics (Table 3), focused on bias, correlation, standard deviation, and normalized mean error. Note that the first three metrics provide independent information about model performance, while normalized mean error contains information about all three previous metrics and is commonly used as a summary metric.

The first group of benchmarks in the PLUMBER experiment were two earlier-generation, physically based models: the Manabe bucket model (Manabe 1969), a simple soil moisture reservoir model with added surface exchange turbulence, and the Penman–Monteith equation (Monteith and Unsworth 1990), which calculates evapotranspiration based on net irradiance, air temperature, humidity, and wind speed. As anticipated (e.g., Chen et al. 1997), modern LSMs outperform these simpler physically based models (Best et al. 2015).

The second group of benchmarks investigated in PLUMBER was those used in PALS (Abramowitz 2012), a web-based database of model simulation and observational land surface datasets, with integrated diagnostic analysis tools. This benchmark group consisted of three empirical models: 1lin, a simple linear regression against downward shortwave radiation; 2lin, a two-dimensional linear regression against downward shortwave radiation and air temperature; and 3km27, a three-dimensional, k-means clustered piecewise-linear regression against downward shortwave radiation, temperature, and relative humidity. All three empirical models were trained and tested with half-hourly flux tower data. Each empirical model was applied out-ofsample separately at each FLUXNET site by calibrating on data from the 19 other sites to establish

	FLUXNET code	Location	Lat	Lon	IGBP land-cover type	Time frame	Years
Amplero	IT-Amp	Italy	41.9041°N	13.6052°E	Croplands	2002-08	4
Blodgett	US-Blo	California, United States	38.8953°N	120.633°W	Evergreen needleleaf	1997-2007	7
Bugac	HU-Bug	Hungary	46.6917°N	19.6017°E	Croplands	2002-08	4
ElSaler2	ES-ES2	Spain	39.2756°N	0.3153°W	Croplands	2004-10	2
ElSaler	ES-ES1	Spain	39.346°N	0.3188°W	Permanent wetlands	1999–2006	8
Espirra	PT-Esp	Portugal	38.6394°N	8.6018°W	Woody savannas	2002-09	4
FortPeck	US-FPe	Montana, United States	48.3077°N	105.102°W	Grasslands	1999–2013	7
Harvard	US-Ha1	Massachusetts, United States	42.5378°N	72.1715°W	Mixed forests	1991-2013	8
Hesse	FR-Hes	France	48.6742°N	7.0656°E	Deciduous broadleaf	1996-2013	6
Howard	AU-How	Australia	12.4943°S	131.152°E	Savannas	2001-13	4
Howlandm	US-Ho1	Maine, United States	45.2041°N	68.7402°W	Mixed forests	1995-2013	9
Hyytiala	FI-Hyy	Finland	61.8474°N	24.2948°E	Evergreen needleleaf	1996-2013	4
Kruger	ZA-Kru	South Africa	25.0197°S	31.4969°E	Savannas	2000-10	2
Loobos	NL-Loo	Netherlands	52.1679°N	5.744°E	Evergreen needleleaf	1996-2013	10
Merbleue	CA-Mer	Ontario, Canada	45.4094°N	75.5187°W	Permanent wetlands	1998-2013	7
Mopane	BW-Ma1	Botswana	19.9165°S	23.5603°E	Savannas	1999-2001	3
Palang	ID-Pag	Indonesia	2.345°S	114.036°E	Evergreen broadleaf	2002-13	2
Sylvania	US-Syv	Michigan, United States	46.242°N	89.3477°W	Mixed forests	2001-09	4
Tumbarumba	AU-Tum	Australia	35.6557°S	148.152°E	Evergreen broadleaf	2000-13	4
UniMich	US-UMB	Michigan, United States	45.5598°N	84.7138°W	Deciduous broadleaf	1998–2013	5

TABLE 1. FLUXNET datasets used in PLUMBER.

regression parameters and then using the meteorological data from the testing site to predict flux variables using these parameters.

The two groups of benchmarks were used to quantify expectations of LSM performance. That is, they provide some understanding of how close to observations we should expect an LSM to be, based on the complexity of the processes at each site and how much information is available in meteorological variables about latent and sensible heat fluxes.

In the PLUMBER experiments, LSMs used the appropriate vegetation type, vegetation height, and reference

height, but otherwise used their default parameter values for the specified vegetation type and selected soil parameter values using their own internal datasets. The LSMs were equilibrated by using the first year of each FLUXNET site repeatedly as a spinup phase. More detail about the PLUMBER experimental protocol can be found in Best et al. (2015).

The results of this comparison are reproduced here for reference in Fig. 2. The columns represent the different LSMs. Within each column, latent and sensible heat fluxes are shown. The vertical axis represents the rank of each LSM for one of these flux

Model	Developer/custodian	Name	Version in PLUMBER
CABLE	Commonwealth Scientific and Industrial Research Organisation (CSIRO)	Community Atmosphere Biosphere Land Exchange model	2.0 and 2.0_SLI
CHTESSEL	European Centre for Medium-Range Weather Forecasts	Carbon Hydrology Tiled ECMWF Scheme of Surface Exchanges over Land	1.1
COLASSiB	Center for Ocean–Land–Atmosphere Studies (COLA)	COLA–Simplified Simple Biosphere (COLA-SSiB)	2.0
ISBA-SURFEX	Centre National de Recherches Météorologiques- Groupe d'Etude de l'Atmosphère Météorologique (CNRM-GAME)	Interactions between Soil, Biosphere, and Atmosphere–Surface Externalisée (ISBA-SURFEX)	3l-7.3 and dif-7.3
JULES	Met Office and Natural Environment Research Council	Joint UK Land Environment Simulator (JULES)	3.1 and 3.1_altP
Mosaic	NASA	Mosaic	1
Noah	Noah	Community Noah land surface model	2.7.1, 3.3, and 3.2
ORCHIDEE	L'Institut Pierre-Simon Laplace (IPSL)	Organizing Carbon and Hydrology in Dynamic Ecosystems (ORCHIDEE)	r1401

TABLE 2. Models used in PLUMBER.

Abbreviation	Formula
MBE	$\frac{\sum\limits_{i=1}^n \left(M_i - O_i\right)}{n}$
NME	$rac{\sum M_i - O_i }{\overline{\Sigma} \overline{O} - O_i }$
sd	$ 1 - rac{\sqrt{\sum M_i - \overline{M}^2}}{\sqrt{rac{\sum O_i - \overline{O}^2}{n-1}}} $
r	$rac{\displaystyle\sum_{i=1}^n (M_i - \overline{M})(O_i - \overline{O})}{\displaystyle\sqrt{\displaystyle\sum_{i=1}^n (M_i - \overline{M})^2}} \sqrt{\displaystyle\sum_{i=1}^n (O_i - \overline{O})^2}$
	Abbreviation MBE NME sd

TABLE 3. Standard statistical set of metrics used in PLUMBER. All metrics are based on half-hourly values. In formulas, *M* represents model data, *O* represents observed flux tower data, and *n* is the number of time steps.

variables, averaged across all four metrics and 20 flux tower sites. Ranks are performed separately for each LSM against the two physically based approaches and the three empirical models, so that the average rank of any of the benchmark models can be different in each LSM. Ranks were used as a way of aggregating performance outcomes across the four metrics and 20 sites.

The key result from PLUMBER, reported by Best et al. (2015), is that LSMs do not perform well in comparison with even simple empirical models for these four common metrics. For sensible heat Q_h , even the simple one-dimensional linear regression against downward shortwave radiation outperforms all of the LSMs (Fig. 2). The slightly more complex 3km27 empirical model outperforms all models for all variables (including net ecosystem exchange of CO₂, not shown here). These results are disturbing, but it is not at all clear from the original experiment what is causing these performance problems, or even if they are particularly meaningful. There are three categories of possible causes of the apparent poor performance by the LSMs:

- the apparent poor performance is due to problems with the PLUMBER methodology;
- the apparent poor performance is due to spurious good performance of the empirical models (e.g., systematic observational error, or empirical models lack of energy conservation constraint); or
- the poor performance is real and is due to poor representations of physical processes, process order, or ability to prescribe appropriate parameter values in LSMs.

Best et al. (2015) did not systematically examine the PLUMBER results in the context of these three categories. Our goal is to either identify the cause of the apparently poor behavior of the LSMs, or—equally usefully—discount possible causes of the problems. Here, we design and execute a number of experiments that target these three categories. As this is a series of discrete experiments, we describe the methods and



FIG. 2. Ranks of LSMs relative to benchmarks, averaged over all metrics and sites [after Fig. 4 in Best et al. (2015)]. Each column shows a different LSM. Within each column, sensible heat (i.e., Q_h) and latent heat (i.e., Q_{le}) are shown. The LSM is in black, and various benchmarks are shown in comparison. The vertical axis shows the average performance rank for each model under four metrics over the 20 FLUXNET site datasets. In each case, a lower value indicates better relative performance. The 3km27 model clearly outperforms the LSMs for both variables, and the two linear regressions consistently outperform all LSMs for sensible heat.



FIG. 3. Histograms of differences between metric values for benchmarks and models with neighboring ranks, for all models at all sites. Values are calculated by taking the difference of the metric value for each model (LSM or one of the five benchmarks) from the model ranked next worst for each LSM, FLUXNET site, metric, and variable. The blue data show the benchmark-to-benchmark metric differences. The red data show the differences between the LSM and the next-worst-ranked benchmark (e.g., if the model is ranked 4, the comparison with the fifth-ranked benchmark). The green data show the difference between the LSM and the next-best-ranked benchmark. Since the models are ordered, all differences are positive (correlation is inverted before differences are calculated).

results together, for each experiment divided into the three categories described above.

2. Methodology and results

a. First possible cause: PLUMBER methodology

There are a number of aspects of the PLUMBER methodology that warrant closer examination. Here we investigate some potentially problematic aspects: the use of ranks instead of metric values, aggregation over sites and metrics, the possibility that PLUMBER was conducted on the wrong time scale, and the simulation initialization procedure.

1) ARE RANKS REPRESENTATIVE?

We first confirm that the PLUMBER ranks are a reasonable representation of the underlying relative real performance values for each metric and variable. PLUMBER used ranks in place of metric values because metric values are not comparable or easily normalizable because of their complex distributions. However, ranks do not necessarily capture all the nuances of the underlying data, and they may misrepresent the performance of the LSMs relative to the benchmarks. For example, if empirical models only outperformed LSMs by very small margins, and when LSMs outperformed empirical models the margins were much larger, the average rank diagnostic could be very misleading.

To assess whether this is a problem in the PLUMBER results, we calculated the differences in metric values between each model (benchmark or LSM) and the next-best and next-worst model. This measure allows us to make statements about the relative performance of the various models, independent of the distribution of the metrics. If, for example, a model appears equally often at each rank, one might expect the distribution of metric margins associated with that model (i.e., "distance" to the next best or worst model) to be similar to the overall distribution of metric margins across all models. This would not be true, however, if the model was consistently only just beating other models, relative to other pairs of models in general. In that case one would expect the distribution of next-worst margins to have a lower mean than overall next-worst distribution, and the distribution of the next-best margin to have a higher mean.

Figure 3 shows the distributions of the differences between each model (benchmark or LSM) and the nextbest and next-worst model. The red and green data highlight the comparisons between the LSMs and the next worst and next best of the five benchmarks, respectively. In general, the red and green have similar distributions, and those distributions are fairly similar to the differences between benchmark pairs (blue histogram), indicating that the ranks are representing the relative performances reasonably well. In cases where the LSM is the worst performing model, there is no red data, and vice versa.

The skew to the right that is clearly visible in most of the plots is to be expected. These metrics all have values that converge on 0 (or 1 in the case of correlation, which is inverted), and become more dense as they approach 0. Therefore, larger differences are to be expected for worse performing pairs of models. Since LSMs tend to perform worse than the benchmarks on average, this 🔶 1lin 🔶 2lin 🔶 3km27 🗢 Manabe_Bucket.2 🔶 model 🗢 Penman_Monteith.1



FIG. 4. As in Fig. 2, but for results where each row represents an individual metric (see Table 3 for metric definitions). The gray line shows the original LSM mean rank for comparison (as in Fig. 2, though note that these data are only comparable with the black line, and not the benchmarks that have also changed).

skew is more pronounced. This suggests that it is unlikely that ranks are unrepresentative of the underlying relative performance differences.

2) IS AGGREGATION OVER SITES AND METRICS PROBLEMATIC?

The results presented in PLUMBER are ranks averaged across multiple metrics and across multiple sites for each variable. It is possible that the averaging process is hiding more distinct patterns of performance—perhaps at particular sites or under particular metrics. To assess whether a particular site or metric was unduly influencing the original PLUMBER results, we reproduce the main PLUMBER plot separately by metric (Fig. 4) and by site (Fig. 5). In both of these plots and in later plots, the original ranks for each LSM from Fig. 2 are shown in gray. Note, however, that the ranks shown in gray are not necessarily ordered with respect to the benchmarks in the same way that they are in Fig. 2 and are only comparable to the black line. For example, in Fig. 2, most LSMs rank better than 2lin for latent heat Q_{le} , but in Fig. 4, the gray line might suggest that some these LSMs performed worse than 2lin, but this is only because the relative rank of 2lin has changed.

Figure 4 shows that while there is some variation between metrics, it is not the case that the LSMs are performing much better or worse than empirical models for any particular metric. Performance relative to the benchmarks is generally mediocre across the board. The



FIG. 5. As in Fig. 2, but for results where each cell represents the average rank of all LSMs at each individual FLUXNET site. The gray line is identical to that shown in Fig. 4.

LSMs do perform better for standard deviation in Q_{le} , outperforming even the 3km27 model in most cases. Best et al. (2015) demonstrated that the LSMs performed better than the empirical benchmarks for the extremes of the distribution of each variable, and our analysis helps confirm that finding. As noted in Best et al. (2015), the empirical models should be expected to produce lower variability since they are regression based. The normalized mean error and correlation metrics were significantly worse than the original aggregate results in Fig. 2. Gupta et al. (2009) showed that RMSE and correlation contain substantially similar information; however, in this study the correlation metric was the least correlated of the four metrics (-0.33 with)mean bias, -0.43 with normalized mean error, and -0.20 with standard deviation difference). On the other hand, correlations between the other three metrics were quite high (0.77 mean bias with normalized mean error, 0.75 mean bias with standard deviation difference, and 0.83 normalized mean error with standard deviation bias). The fact that the LSMs appear to be performing best under two of these three highly correlated metrics (mean bias and standard deviation difference), at least relative to the 3km27 benchmark, may indicate that the PLUMBER results overestimate LSM performance.

Figure 5 shows that there is considerable diversity of performance between sites for the LSMs. In this case, results are averaged over all 13 LSMs and the four metrics in Table 3. For example, the LSMs perform relatively very well for Q_h at the ElSaler site. This site is unusual: it is situated on a low-lying narrow spit of land between a small lake and the Mediterranean Sea and is likely heavily influenced by horizontal advection. It is possible that rather than the LSMs performing well here, it is actually the empirical models that are performing poorly because they were calibrated on all other sites that do not exhibit behaviors seen at ElSaler. This possibility is supported by the fact that the models that include some measure of humidity (3km27 and Penman-Monteith) perform worse than the simpler linear regressions. ElSaler2 is another unusual case, an irrigated cropland site in Mediterranean Spain. The LSMs and Manabe bucket model, which do not have information about the additional water input to the system, do very poorly. The unconstrained reservoir in the Penman-Monteith equation in this case works very well. There are a number of sites where LSMs consistently perform poorly-Espirra provides an example pattern that we might expect from the original PLUMBER results—with LSMs performing worse than empirical models, but much better than early theoretical models. However, there are other sites where LSMs are performing poorly even against the older approaches, especially for Q_h , such as Amplero and Sylvania, and there are no sites where LSMs perform consistently well relative to the benchmarks for both fluxes. While each of these breakdowns—by metric and by site—give us some insight into how LSMs are behaving, they do not explain the cause of the general pattern of apparent poor performance.

3) Do LSMs perform better on longer time scales?

Another possibility is that poor performance in the short-time-scale, half-hourly responses of LSMs are dominating the performance metrics. While versions of these models are designed for both climate and weather prediction, here we are largely concerned with longterm changes in climate and the land surface. In this context, short-time-scale responses may be relatively inconsequential, as long as the longer-term result is adequate. It is plausible, for example, that short time lags in various state variables built into LSMs might be adversely affecting the half-hourly model performance, while improving the longer-time-scale skill of the model. All of the original PLUMBER metrics are calculated on a per time step basis, and so do not take this possibility into account. To examine this, we recalculate the PLUMBER ranks after first averaging the half-hourly data to daily, monthly, and seasonal time steps.

Figure 6 reproduces the PLUMBER plots after averaging data to three different time scales: daily averages, monthly averages, and seasonal averages. While there are some changes in these plots, there is no major improvement of LSM behavior relative to the empirical benchmarks. On all time scales, the LSMs are consistently outperformed by the empirical benchmarks, suggesting that the problems found in PLUMBER are not related to time scale.

4) ARE INITIAL CONDITIONS A PROBLEM?

It is possible that the initialization procedure used in PLUMBER is inadequate. If the spinup period was not long enough for state equilibration, or it was not representative of the period immediately preceding the simulation, then we would expect to see a stronger bias in the early parts of the first year of the data for each run. PLUMBER used a spinup procedure that involved repeatedly simulating the first year at each site 10 times before running over the whole period and reporting model output. To test whether poor spinup might be the cause of the poor performance seen in PLUMBER, we calculated a number of new metrics over each simulation, for each variable, based on daily average data. First, we calculate the day at which each of these simulation time series first crosses the equivalent observed

🔶 1lin 🔶 2lin 🔶 3km27 🗢 Manabe_Bucket.2 🔶 model 🗢 Penman_Monteith.1



FIG. 6. As in Fig. 2, but for values that are averaged over daily, monthly, and seasonal time periods. The gray line is identical to that shown in Fig. 4.

time series, both as an absolute value and as a percentage of the length of the dataset, which gives some indication of whether the simulation has converged on the observed data. Next, we calculate the difference in slope parameters of a linear regression over the two time series, and also the significance of this difference (where the null hypothesis is no difference). Last, we check if the bias is decreasing—that is, if the simulations have positive mean errors, is the trend slope negative (e.g., mean error is closer to zero in the second half of the time series) or vice versa?

Figure 7 shows the results of the approaches described above. For each of the two fluxes (rows), using daily average data, it shows the first day in the time series that the simulated flux is equal to, or crosses, the observed flux (first column, logarithmic scale) and the results expressed as a percentage of the time series (second column); the difference in the slopes of linear regressions of simulated and observed series over time (third column; $W day^{-1}$); significance of the difference in the previous metric (fourth column; values left of the red line are significant at the $\alpha = 0.05$ level; ~44% of all values); and the rate at which the bias is decreasing, measured by means of model error divided by the gradient of model error (fifth column; negative values indicate the simulations have a trend toward the observations). Each panel is a histogram, with each entry colored by the FLUXNET site it represents.

The first two metrics show that in nearly all cases, the simulations' time series quickly cross the observed time series (76% of simulations cross in the first 1% of the period, and 97% cross in the first 10%), indicating that it is unlikely that lack of equilibration explains the poor behavior of the LSMs relative to the benchmarks. The



FIG. 7. Histograms of model spinup metrics, based on daily averages, from all LSMs at all sites. From left to right: day at which the simulated series crosses the observed series; day at which the simulated series crosses the observed series, but as a percentage of the time series; difference in the slopes of linear regressions of simulated and observed series over time (W day⁻¹); significance of the difference in the previous metric (values left of the red line are significant at the $\alpha = 0.05$ level; ~44% of all values); and the rate at which the bias is decreasing, measured by mean(error)/slope(error) (negative values indicate the simulations have a trend toward the observations). Colors indicate the FLUXNET site at which the simulation is run.

TABLE 4. Correlation between model metrics in Fig. 7.

	First crossing	First cross percent	Slope diff	Slope diff significance
Bias decreasing	-0.017	-0.019	-0.025	-0.006
First crossing		0.990	0.029	0.0386
First cross percent			0.015	0.031
Slope diff				0.034

third and fourth metrics show the differences between the trends in the observations and the simulations and the significance of those differences. In the majority of cases, effect sizes are quite small, with 61% of absolute trend differences less than $0.02 \text{ W} \text{ day}^{-1}$ or $7.3 \text{ W} \text{ yr}^{-1}$ (third column, Fig. 7), which is well within the standard error of the time series. Forty-five percent of these trend differences are significant at the $\alpha = 0.05$ level (fourth column, Fig. 7), but there is no indication of a pattern of trends toward a lower bias; 54% of simulations have a trend that increases rather than decreases the bias (column 5). The colors in the plot specify the FLUXNET sites, and as indicated, aside from the two first-crossing metrics, there is very low correlation between metrics ($r \ll 0.05$, see Table 4).

We have therefore not been able to find obvious major systematic flaws in the PLUMBER methodology. The poor performance of the LSMs in PLUMBER, relative to the empirical benchmarks, cannot be dismissed based on any obvious flaw in the methodology.

b. Second possible cause: Spurious empirical model performance

We next examine the possibility of spurious good performance by the empirical models. While there are a number of possibilities related to data quality, we focus on one main possibility that has been brought up multiple times by the community in response to the original PLUMBER paper.

LACK OF ENERGY CONSERVATION CONSTRAINTS

The obvious candidate is that the empirical models are able to perform so well relative to the LSMs because they do not have any kind of built-in constraint for energy conservation. This allows them to potentially produce results that predict individual flux variables quite well, but are physically inconsistent (e.g., outgoing flux energy is not constrained by net radiation). One way to test this hypothesis is to build empirical models that have additional constraints that ensure that energy is conserved.

Because of the effects of energy storage pools (mainly in the soil), it is not a trivial matter to produce a conservation-constrained empirical model. We therefore approach the problem from the opposite direction: we assume that energy conservation in the LSMs is correct and use the calculated available energy $(Q_h + Q_{le})$ from each LSM to constrain the empirical model output:

$$Q_{\rm emp}' = \frac{Q_{\rm emp}}{(Q_{h_{\rm emp}} + Q_{\rm le_{\rm emp}})} (Q_{h_{\rm LSM}} + Q_{\rm le_{\rm LSM}}),$$

where Q_{emp} can be either $Q_{h_{emp}}$ or $Q_{le_{emp}}$. An alternative approach might be to correct the observations with the LSMs' total energy and to retrain the empirical models on the corrected data. We have no a priori reason to expect that this approach would provide qualitatively different results, and it would require significantly more computation.

Our approach effectively forces each empirical model to have the same radiation scheme and ground heat flux as the LSM it is being compared to (since available energy $Q_{le} + Q_h$ is now identical) and preserves only the Bowen ratio from the original empirical model prediction. While this makes the empirical models much more like the LSMs, it informs us whether the empirical models were simply reproducing a systematic lack of energy conservation in the flux tower data. That is, if these modified empirical models perform similarly to their original counterparts, then energy conservation, while no doubt a real data issue, is not the cause of this result. If the reverse is true-that the modified empirical models no longer outperform the LSMs-there are at least two possibilities. Most obviously, the empirical models may indeed be fitting to systematically biased observational data. Alternatively, poor available energy calculations on the part of LSMs might cause the degradation of the modified empirical models, so that energy conservation is less of an issue. There are some difficulties with the transformation shown in the equation above. When the denominator in this equation approaches zero the conversion could become numerically unstable. Under these conditions, we replace all values of Q_h and Q_{le} with the values from the LSM whenever $|Q_{h_{\text{emp}}} + Q_{\text{le}_{\text{emp}}}| < 5 \,\text{W}\,\text{m}^{-2}$. This effectively means that only daytime values are modified.

If the energy-conserving empirical models still outperform LSMs, it would indicate that calculation of available energy in LSMs is relatively sound and that the

🔶 1lin 🔶 2lin 🔶 3km27 🔶 Manabe_Bucket.2 🜩 model 🔶 Penman_Monteith.1



FIG. 8. As in Fig. 2, but for energy conservation constrained empirical models. The gray line is identical to that shown in Fig. 4.

energy partitioning approach is the likely cause of the poor performance. That is, even when empirical models are forced to have the same available energy as each LSM, performance ranks are essentially unchanged. Alternatively, if the energy-conserving empirical models perform poorly, it may either indicate that empirical models are trained to match systematically biased, nonconserving flux tower data or that the calculation of available energy in LSMs is the main cause of their poor performance.

The results of the energy-conserving empirical model experiment are shown in Fig. 8. We wish to reinforce that Fig. 8 shows precisely the same LSM, Manabe bucket, and Penman–Monteith simulations as Fig. 2, and only the empirical benchmarks have changed (which in turn affects the other models' ranks).

It is clear that this change to the empirical models offers some LSMs a relative improvement in their rank. Noah2.7.1 and ORCHIDEE now beat all empirical models for Q_{le} , for example. This is far from a uniform result, however. Note also that Q_{le} performance from CABLE2.0_SLI, ISBA-SURFEX31, and Noah3.2 is now worse than 2lin, which was not the case in Fig. 2. The energy constraint has actually improved the empirical model performance in these cases. It is also still the case that all LSMs are outperformed by the energy-conserving versions of 1lin for Q_h . It therefore appears unlikely that the energy conservation issues in flux tower data are the cause of the empirical models' good performance.

While some of the changes seen in Fig. 8 can be attributed to the forcing of energy conservation on empirical models, there are other possible interpretations. They could be reflecting the effect that each LSM's available energy calculation had on the empirical models. For example, if a particular LSM had a very poor estimate of instantaneous available energy (i.e., $Q_{le} + Q_h$) because of issues in its radiation or soil heat transfer schemes, forcing this estimate on all of the empirical models might degrade their performance in a nonphysical way. This would of course appear in Fig. 8 as a relative improvement in the LSM's performance. It is not clear whether this, or accounting for a lack of energy conservation in empirical models, is the cause of the improvements and degradations in performance we see in Fig. 8.

One unavoidable problem with this methodology is that if the flux tower data have a consistent bias in the evaporative fraction, then the LSMs will appear to perform relatively worse because of the empirical models overfitting that bias. Figure 9 shows the biases in simulated evaporative fraction at each site across all LSMs. This plot consists of standard box plots showing the mean, first and third quartiles, and outliers. The biases are calculated by taking

$$\left(\frac{Q_{\mathrm{le_{sim}}}}{Q_{h_{\mathrm{sim}}}+Q_{\mathrm{le_{sim}}}}-\frac{Q_{\mathrm{le_{obs}}}}{Q_{h_{obs}}+Q_{\mathrm{le_{obs}}}}\right)$$

using daily data and excluding all cases where $|Q_h + Q_{le}| < 1 \text{ W m}^{-2}$ for either simulations or observations, to avoid numerical instability. It is clear that at some sites the LSMs have an apparent bias in evaporative fraction. It is not possible to be certain whether this bias is in the flux tower data or because of shared problems between the LSMs. We address this in the discussion.

This analysis indicates that, while problems with the flux tower data may contribute in a small way, they do not explain the entirety of the poor performance seen in PLUMBER. In general, the LSMs are not only predicting total heat poorly, they are also predicting the partitioning of that heat poorly.

c. Third possible cause: Poor model performance

Finally, we search for indications that the problem might lie with the LSM simulations themselves. We examine two possibilities: LSM performance over short time scales and performance at different times of the day. We also explore how the LSMs perform as an ensemble, in an attempt to assess whether problems might be shared across models.

1) How do LSMs perform over short time scales?

When investigating the PLUMBER methodology, as outlined above, we examine whether short-time-scale



FIG. 9. Biases in daily evaporative fraction for each LSM simulation, grouped by site.

variability is dominating the PLUMBER metrics by averaging data to different time scales before recalculating performance measures. The inverse of this possibility is that rather than getting the short-time-scale aspects of climate wrong, the LSMs are actually simulating the highfrequency responses well, but failing over the long term. This would occur, for example, if the magnitude of the soil moisture reservoir were the wrong size, or the input or output to this reservoir caused it to dry too quickly or too slowly. To test this possibility, we remove all of the low-frequency variability from the error time series, by first bias correcting the simulation on a daily basis for each variable $(Q'_{sim} = Q_{sim} - \overline{Q}_{sim} + \overline{Q}_{obs})$, for each day) and then removing the average daily cycle over the remaining residuals. This gives us a model time series that has the same mean daily temperature and average daily cycle as the observations, but retains all of the modeled high-frequency variability.

The high-frequency-only results are shown for each metric in Fig. 10. Because of the nature of the bias correction, the bias metric (second row in Fig. 4) is always 0 for the LSMs, resulting in a trivial rank of 1, and so we remove the bias metric from these results. The effect this has can be seen by comparing Fig. 10 to the first, third, and fourth rows in Fig. 4. In all three metrics there are notable improvements in LSM ranks (averaged over all sites), suggesting that a significant portion of LSM error is likely due to the modulation of instantaneous model responses by the model states (e.g., soil moisture and temperature). The degree of improvement does vary between models to some degree—CABLE2.0_SLI,

COLASSiB, and Noah3.3 improved in absolute rank in all metrics as a result.

2) DO LSMS PERFORM BETTER AT DIFFERENT TIMES OF THE DAY?

The LSMs appear to be having problems extracting all of the available information from the available meteorological forcings, especially downward shortwave radiation (SW_{down}), as evidenced by the 1lin model outperforming each LSM for Q_h . It thus seems likely that the LSM performance might vary according to the availability of that information. To test this possibility, we split the analysis over time of day, splitting each time series into night (2100–0300 LT), dawn (0300–0900 LT), day (0900–1500 LT), and dusk (1500–2100 LT) and repeating the analysis for each subseries.

The time-of-day analysis is presented in Fig. 11. As might be expected, there is clear variation in LSM performance relative to the benchmarks at different times of the day. The LSMs generally outperform the 1lin and 2lin models at nighttime. This is to be expected, as these two benchmarks, 1lin especially, have essentially no information at this time of day. In general, the LSMs all appear to be having difficulty with both fluxes around sunrise. It is worrying that some of the LSMs appear to be doing worse than a linear regression on sunlight during the nighttime for latent heat (COLASSiB, ISBA-SURFEX31, and ORCHIDEE). However, the performance differences are small in those cases and may be simply an artifact of the data (e.g., the empirical models fitting noise in FLUXNET).

🔶 1lin 🔶 2lin 🔶 3km27 🔶 Manabe_Bucket.2 🔶 model 🔶 Penman_Monteith.1



FIG. 10. As in Fig. 2, but for high-frequency response only, by metric; LSMs are bias corrected on a daily basis and then have the daily cycle in the errors removed. The gray line is identical to that shown in Fig. 4. The mean bias error metric is not included because it is trivially zero because of the bias correction process.

Overall, it does not appear to be the case that the LSMs are performing well at any particular times of the day.

3) How do the LSMs perform as an ensemble?

Last, we investigate whether the nature of the poor performance is a problem that is shared among models by examining the performance of the LSMs as an ensemble. Model ensemble analysis has a long history in the climate sciences (e.g., the Coupled Model Intercomparison Project; Meehl et al. 2007; Taylor et al. 2012), as well as in the land surface modeling community (Dirmeyer et al. 2006). Ensemble analysis allows us to identify similarities in performance between the LSMs. If each LSM is performing poorly for very different reasons, we might expect that at a given site, the time series of model error (model observed) between different models would be uncorrelated. If this were the case, the multimodel mean should provide a significantly better estimate of the observed time series, since the eccentricities causing each model's poor performance will tend to cancel each other. By analogy, the standard deviation of the mean of *n* random number time series, each with standard deviation 1 and mean 0, is $1/\sqrt{n}$. As an attempt to try to ascertain the degree of shared bias among LSMs, we choose to



FIG. 11. As in Fig. 2, but for results split by daily cycle. The four rows represent the 6-h periods around dawn (0300–0900 LT), noon (0900–1500 LT), dusk (1500–2100 LT), and midnight (2100–0300 LT). The gray line is identical to that shown in Fig. 4.

examine three different ensemble means: the unweighted average; the error-variance-based, performance-weighted mean; and the error-covariance independence-weighted mean (Bishop and Abramowitz 2013; Haughton et al. 2015). A priori, we should expect these ensemble means to perform differently in different circumstances. First, as mentioned above, if errors from different models have pairwise low correlations, we should expect the model mean to perform better than individual models. Next, if there are substantial differences in performance of the models, we should expect the performance-weighted mean to outperform the unweighted mean. If performance across the ensemble is similar but errors are highly correlated in a subset of the LSMs, then we should expect the independence-weighted mean to outperform both the unweighted mean and performance-weighted mean. The corollary is that if the independence-weighted mean does not outperform the unweighted mean, this likely indicates that problems causing poor performance are shared among LSMs.

The results of the performance of the three ensemble means are shown in Fig. 12. The means all perform similarly, or slightly better than the best LSMs under each metric (see Fig. 4). However, the means are still outperformed by the empirical models in many cases. It is notable that there is also very little improvement under either of the weighted means. The performance-weighted mean only gives a slight improvement, which confirms that the differences in performance between LSMs relative to the benchmarks are not significant. The independence-weighted mean also has little improvement, which gives an indication that problems with performance are shared across LSMs.

3. Discussion

The PLUMBER results are worrisome, and it seems sensible to approach them with some skepticism. It is tempting to write off the results as an artifact of the PLUMBER methodology, but this does not appear to be the case. Over all LSMs tested, there is a consistent problem of poor performance relative to basic empirical models that is not obviously related to simulation initialization, particular sites or metrics biasing the analysis, or the time scale of the analysis. Despite the very wide range of performance ranks across different flux tower sites, once the obvious, understandable cases are removed (especially the ElSaler, ElSaler2 pair of sites, for different reasons), the aggregated picture of performance in Fig. 2 seems broadly representative of our current LSMs.

In our energy-conserving empirical model analysis, we rescaled the total available energy in the empirical models to match that in each LSM, effectively making the total available energy identical in each pair of models and only comparing the partitioning of that energy into Q_h and Q_{le} . We then showed that there are biases between the LSMs and the FLUXNET data, but that across sites there is no consistent bias that might cause the empirical models to perform spuriously well. There are known problems with energy conservation in flux tower data- $R_{\text{net}} = Q_{\text{le}} + Q_h + Q_g$ is unbalanced by 10%–20% at most sites (Wilson et al. 2002). However, this does not tell us anything about any potential bias in the evaporative fraction. Indeed, Wilson et al. (2002) note that the flux biases are independent of the Bowen ratio. Other studies have found that energy balance closure is dependent on stability (Kessomkiat et al. 2013; Stoy et al. 2013). We corrected the empirical model with the evaporative fraction, which is very close but more stable than the Bowen ratio suggested by Wilson et al. (2002). There is, however, discussion in the literature that eddy flux measurements might underestimate sensible heat much more than latent heat (e.g., Ingwersen et al. 2011; Charuchittipan et al. 2014; Mauder and Foken 2006). This would affect the PLUMBER results for sensible heat and might improve LSM ranks. It would not affect the latent heat results, however, and LSMs would still perform worse than the empirical benchmarks for the normalized mean error and correlation metrics.

So, if there is a problem with the LSMs, as appears to be the case, where does it leave us? There are two broad possibilities to investigate.

The first, and perhaps most confronting, is that there are flaws in the structuring, conception of the physics, or ordering of processes in the models. The results from the three approaches to LSM averaging suggest that such a problem might be largely shared among LSMs. LSMs do commonly share some similar conceptualizations of land surface processes, even if they do not share implementation details. Masson and Knutti (2011) showed how interrelated climate models can be. Those results include many of the models used here, and it would be interesting to see such an analysis performed on LSMs alone.

Examples of such shared problems might be that all of the LSMs could be missing a major component or a relationship between components, or they may share a flawed representation of one or more components. This part of the modeling process is hard to analyze rigorously; however, some analysis of assumptions contained in models and the effects that those assumptions have on model performance has been undertaken (e.g., Clark et al. 2008; De Kauwe et al. 2013; Zaehle et al. 2014). In



1lin 🔶 2lin 🔶 3km27 🔶 Manabe_Bucket.2 🔶 model 🔶 Penman_Monteith.1

FIG. 12. As in Fig. 2, but for the results for three different means across all LSMs, by metric. The gray line is identical to that shown in Fig. 4. In general, we should expect means to perform better under all metrics except the standard deviation metric, as the averaging process acts as a smoother, removing noncorrelated noise from the model results.

principle, one could take a single LSM and replace major model components with calibrated linear regressions (if the observational data were available to create these) and compare performance, in order to pinpoint which component is the main cause of the poor performance. This would likely require a quantity of process-level data that is not yet available.

While we largely present negative results in our attempts to pinpoint these problems, there are some indications as to where the problem may lie if model physics is the cause of this result. The energy-conserving empirical models give a strong indication that the calculation of available energy for Q_{le} and Q_h is not the main problem. That is, since the conserving empirical models effectively have the same R_{net} and ground heat flux as the LSMs and still broadly outperform the LSMs, we assume that the main issue is in the calculation of these fluxes. While there are snow periods in some of these datasets, the majority does not include any significant snow-we can probably safely ignore snow submodels as a cause of the overall result. It does appear that there are some issues in the available energy calculations that vary across models. Some models, for example, do perform better in a relative sense once the empirical models are forced to match their available energy (cf. Figs. 2 and 8). Overall, however, this does not make a qualitative difference to LSM ranks against the empirical models. The analysis removing diurnal means (Fig. 10) also broadly supports the idea that available energy and partitioning is being adversely affected by storage. That is, when the error in the diurnal average and average diurnal cycle was removed from LSMs, effectively removing any bias from inappropriate soil moisture levels and leaving behind only each LSM's high-frequency responses, there was an improvement in performance. Ideally, we would like to test directly whether, for example, soil moisture is correlated with the accuracy of evaporative fraction prediction. Unfortunately, the FLUXNET datasets we used did not all contain soil moisture observations. In the cases that did report soil moisture, major challenges exist in using these data to evaluate LSMs. Observations are taken over different depths, using different measurement strategies, for example. There are also major issues in what soil moisture means in an LSM (Koster et al. 2009) and whether this variable can be compared directly with observed soil moisture. We therefore avoid comparisons of the LSM results with observed soil moisture but note that if the problems of data quality, consistency of measurements, and issues of scale can be resolved, this would provide a particularly good way forward for resolving why the LSMs perform poorly.

One caveat that must be added here is that these simulations are all run off-line, uncoupled from an atmosphere model. In climate simulation and numerical weather prediction experiments, the LSM would be coupled to an atmosphere model that provides feedback to the land surface in a way that fixed meteorological forcings cannot, and this feedback may provide damping of errors that the LSMs produce. Wei et al. (2010) indicates an effect along these lines in dry regions, by showing that an ensemble of LSMs coupled to an atmosphere model can produce higher variance between the LSMs when they are coupled individually, likely due to the fact that the strength of the coupling feedback is divided among the participating LSMs. Holtslag et al. (2007) also find that coupled models tend to produce less variance in stable boundary layer conditions because the fluctuating surface temperature provides feedback to the heat fluxes. A logical next step is therefore to perform a PLUMBER-like benchmarking evaluation in a coupled environment. Because of the difficulty of coupling many LSMs with one or more atmosphere models, as well as the problem of how to fit the benchmarks, such an experiment would be extremely challenging to undertake.

Calibration is also an ongoing problem, particularly because of the large number of poorly constrained parameters and internal variables, combined with the nonlinearity of the models, which leads to problems of equifinality. These results might also reflect the compensating effect of calibration against streamflow or gridded evapotranspiration products, where model structural and spatial property assumptions form part of the calibration process. Experiment-specific calibration may have improved the performance of the LSMs in PLUMBER. However, calibrating LSMs per site would give them an unfair advantage over the empirical models, which are only calibrated out of sample and which use no site-characteristic data. The simulations in PLUMBER were run with appropriate reference heights and IGBP vegetation type, using the LSM's default calibration for that vegetation type. Soil characteristics were selected by individual modeling groups. Clearly, using broad vegetation classes risk losing a lot of site-level specificity, but there is no way to calibrate the LSMs for specific sites while ensuring no overfitting (e.g., out-of-sample calibration) within the PLUMBER dataset, since there are not multiples of each vegetation class represented. Improved per-vegetation class calibration using other FLUXNET sites may help, but at least some of the LSMs in this study are already calibrated on FLUXNET or similar datasets at multiple sites and should perform reasonably well over these 20 datasets without recalibration. While there are advanced methods of multicriteria calibration available (e.g., Guerrero et al. 2013; Gupta et al. 1999), as well as viable alternatives to performance-based calibration (Schymanski et al. 2007), it would seem sensible to also focus on model parsimony, especially in components that are largely underconstrained. However, even if calibration is part of the problem here, it must be remembered that the empirical models are acting on only 1-3 of the 7 meteorological variables available to the LSMs, and also take no account of spatial or temporal variables. While it is true that adding further forcing variables would not guarantee a better result, for example, if those variables have systematic errors, the consistency of performance of the empirical models indicates that that is not the case for at least downward shortwave radiation, air temperature, and relative humidity, and we have no a priori reason to expect it to be the case with the other variables.

It is also worth reflecting on the fact that the core conceptual process representations in LSMs were derived before any high-density data were widely available across different biomes. While the majority of these LSMs are calibrated on some site-level data, there is the possibility that our conceptually consistent LSMs are in some way not physically consistent with observations. An example of this possibility, that may explain the PLUMBER result that the LSMs are almost always worse at simulating Q_h compared to Q_{le} , relates to how the models are designed. The formulation of Q_h and Q_{le} in LSMs commonly refers to a "within canopy temperature," for example, through which these fluxes are exchanged with the atmosphere above the canopy. Imagine that this within canopy air temperature is erroneous. Under these circumstances, Q_h would systematically be simulated poorly relative to Q_{le} , because it is not limited by available moisture. On top of this, energyconservation correction formulas may be partitioning the conservation error poorly.

We cannot test this in all models involved in PLUMBER, but we can test this idea using one of the PLUMBER models. We took CABLE and introduced an error in the initial temperature of the canopy air space ranging from -5 to +5 K, at the start of each time step, and we then examined the impact of this error on Q_h and Q_{le} . Figure 13 shows how the error in Q_h and Q_{le} scales with the error in within canopy air temperature and shows that the error in Q_h increases much more quickly than the error in Q_{le} . We are not suggesting here that this is why all LSMs testing in PLUMBER show this behavior, but we do suggest that there are key variables, common to LSMs, that act as pivots in the performance of an LSM and that are not resolved by feedbacks. While canopy interception cannot introduce too large an error (because too much evaporation in one hour will be compensated by too little in the next hour), if a systematic error is implicit in the interpolation of a reference air temperature to a canopy air temperature, then this may not be compensated by feedbacks and lead to an error that is not resolved on longer time scales. We can demonstrate this for CABLE, and we suggest it is a plausible explanation for other LSMs. We suspect that other similar pivot variables, not ameliorated by feedbacks, might exist and might provide keys to unlocking the PLUMBER results.

The second possibility is that the LSMs are conceptually correct but are too complex for the task at hand. Modern LSMs have around 40 spatially varying parameters. At the scales that they normally operateglobally or regionally-observations rarely adequately constrain these parameters. To get around this issue they are usually calibrated, often using flux tower data, for each vegetation type. This process makes assumptions about landscape homogeneity and forces the LSM to behave consistently with the time, place, and circumstances of the calibration data. Using complex LSMs in this way may be forcing relatively capable models to operate essentially as empirical models, and using them out of sample. If we only use very simple metrics this can appear to be an issue of equifinality in calibration, but in reality the right answer is obtained for the wrong reasons, and as a result poor predictive outcomes are likely.

If true, this suggests that the appropriate level of complexity for a global LSM is a model with a parameter set of approximately the same dimension as the number of independent observable surface properties at the global scale—perhaps an order of magnitude smaller than modern LSMs today. While this is approximately the amount of information we provide LSMs at this scale, by prescribing vegetation and soil types, it is the fixed parameters, or forced covariation of these parameters, that is potentially more important. Related issues of poor parameter constraint were explored by Mendoza et al. (2015). It should also be noted that regression methods, which are based on maximizing variance of the variables we attempt to predict, benefit from a simpler method of fitting and can make stronger use of some observed variables that are not pure predictors, such as relative humidity, which is highly correlated with the Bowen ratio (Barros and Hwu 2002), and therefore may have a substantial advantage. However, this only explains the performance of the 3km27 benchmark and not the fact that the simpler regressions still outperform the LSMs for Q_h .

It is also possible that the problems identified by PLUMBER do not have a single cause and are simply an agglomeration of small, individually insignificant errors, including some of those possibilities identified here. While our results do not explicitly resolve the performance problems shown in the original PLUMBER results, they do help us to rule out a number of possible causes, and in doing so, suggest directions for further investigation.

4. Conclusions

We investigated three broad categories of possible causes for the key result in the original PLUMBER experiment-LSMs being outperformed by simple, outof-sample empirical models. These were the experimental methodology of PLUMBER; spurious good performance of the empirical models in PLUMBER resulting from systematic bias in flux tower data; and genuine poor performance of LSMs. While not every aspect of PLUMBER methodology was investigated, we did establish that particular sites or metrics were not biasing the result. Analyzing data on different time scales similarly had little effect, and there did not appear to be any systematic drift toward observed values that might be indicative of a systematic failure in the model spinup protocol. We also repeated the experiment with energy-conserving versions of the original empirical models used in PLUMBER, constrained by the available energy calculations of each LSM, to try to ascertain whether a lack of energy conservation on the part of empirical models was the likely cause. Again, this had little effect on the result.

This leaves only the last of these three causes, the LSMs themselves. The empirical models suggest that there is more information in the input data available to reproduce observed latent and sensible heat than the LSMs are using. The calculations of the heat fluxes and



FIG. 13. Mean error in Q_h and Q_{le} as a result of perturbing the initial canopy air temperature at each time step, in CABLE at the Tumbarumba site in southeastern Australia. Temperature was perturbed by $\pm(5, 2, 1, 0.5, 0.2)$ K, and a control run is included. All model parameters were left as default values. The response in Q_h to negative temperature perturbations is about 50% stronger than in Q_{le} , and about 3 times stronger for positive perturbations.

the model states upon which these depend are therefore the most likely candidates for the cause of the large discrepancies observed here. It remains a topic for further investigation whether this is ultimately the result of, for example, overparameterization, missing process, problems with calibration, or one of several other possible reasons. Not all models are developed with the same purpose, and some LSM development may have focused on very different aspects of the model, such as the distribution of natural vegetation, which might lead to models that are conceptually consistent but observationally inconsistent when predicting heat fluxes. We cannot recommend specific LSM improvements, but rather provide a framework for model developers against which they can check their developments.

The validity of the benchmarking methodology in Best et al. (2015) was further evaluated in this study. It is worth noting that while PLUMBER may have undiscovered flaws, it is still extremely valuable: the relative poor performance of LSMs would likely have remained hidden under any previous model evaluation or intercomparison methodology.

Acknowledgments. We acknowledge the support of the Australian Research Council Centre of Excellence for Climate System Science (CE110001028). M. Best and H. Johnson were supported by the Joint DECC/ Defra Met Office Hadley Centre Climate Programme

(CA01101). This work used eddy covariance data acquired by the FLUXNET community and in particular by the following networks: AmeriFlux [U.S. Department of Energy, Biological and Environmental Research, Terrestrial Carbon Program (DE-FG02-04ER63917 and DE-FG02-04ER63911)], AfriFlux, AsiaFlux, CarboAfrica, CarboEuropeIP, CarboItaly, CarboMont, ChinaFlux, FLUXNET-Canada (supported by CFCAS, NSERC, BIOCAP, Environment and Climate Change Canada, and NRCan), Green-Grass, KoFlux, LBA, NECC, OzFlux, TCOS-Siberia, and USCCC. We acknowledge the financial support to the eddy covariance data harmonization provided by CarboEuropeIP, FAO-GTOS-TCO, iLEAPS, Max Planck Institute for Biogeochemistry, National Science Foundation, University of Tuscia, Université Laval and Environment and Climate Change Canada, and U.S. Department of Energy and the database development and technical support from Berkeley Water Center; Lawrence Berkeley National Laboratory; Microsoft Research eScience; Oak Ridge National Laboratory; University of California, Berkeley; and University of Virginia.

REFERENCES

Abramowitz, G., 2012: Towards a public, standardized, diagnostic benchmarking system for land surface models. *Geosci. Model Dev.*, 5, 819–827, doi:10.5194/gmd-5-819-2012.

- Barros, A. P., and W. Hwu, 2002: A study of land-atmosphere interactions during summertime rainfall using a mesoscale model. J. Geophys. Res., 107, 4227, doi:10.1029/ 2000JD000254.
- Best, M. J., and Coauthors, 2015: The plumbing of land surface models: Benchmarking model performance. J. Hydrometeor., 16, 1425–1442, doi:10.1175/JHM-D-14-0158.1.
- Bishop, C. H., and G. Abramowitz, 2013: Climate model dependence and the replicate Earth paradigm. *Climate Dyn.*, 41, 885–900, doi:10.1007/s00382-012-1610-y.
- Charuchittipan, D., W. Babel, M. Mauder, J.-P. Leps, and T. Foken, 2014: Extension of the averaging time in eddycovariance measurements and its effect on the energy balance closure. *Bound.-Layer Meteor.*, **152**, 303–327, doi:10.1007/ s10546-014-9922-6.
- Chen, T. H., and Coauthors, 1997: Cabauw experimental results from the Project for Intercomparison of Land-Surface Parameterization Schemes. J. Climate, 10, 1194–1215, doi:10.1175/ 1520-0442(1997)010<1194:CERFTP>2.0.CO;2.
- Clark, M. P., A. G. Slater, D. E. Rupp, R. A. Woods, J. A. Vrugt, H. V. Gupta, T. Wagener, and L. E. Hay, 2008: Framework for Understanding Structural Errors (FUSE): A modular framework to diagnose differences between hydrological models. *Water Resour. Res.*, 44, W00B02, doi:10.1029/ 2007WR006735.
- De Kauwe, M. G., and Coauthors, 2013: Forest water use and water use efficiency at elevated CO₂: A model–data intercomparison at two contrasting temperate forest FACE sites. *Global Change Biol.*, **19**, 1759–1779, doi:10.1111/gcb.12164.
- Dirmeyer, P. A., 2011: A history and review of the Global Soil Wetness Project (GSWP). J. Hydrometeor., 12, 729–749, doi:10.1175/JHM-D-10-05010.1.
- —, A. J. Dolman, and N. Sato, 1999: The pilot phase of the Global Soil Wetness Project. *Bull. Amer. Meteor. Soc.*, **80**, 851–878, doi:10.1175/1520-0477(1999)080<0851:TPPOTG>2.0.CO;2.
- —, X. Gao, M. Zhao, Z. Guo, T. Oki, and N. Hanasaki, 2006: GSWP-2: Multimodel analysis and implications for our perception of the land surface. *Bull. Amer. Meteor. Soc.*, 87, 1381– 1397, doi:10.1175/BAMS-87-10-1381.
- Guerrero, J.-L., I. K. Westerberg, S. Halldin, L.-C. Lundin, and C.-Y. Xu, 2013: Exploring the hydrological robustness of model-parameter values with alpha shapes. *Water Resour. Res.*, 49, 6700–6715, doi:10.1002/wrcr.20533.
- Guo, Z., and Coauthors, 2006: GLACE: The Global Land– Atmosphere Coupling Experiment. Part II: Analysis. J. Hydrometeor., 7, 611–625, doi:10.1175/JHM511.1.
- Gupta, H. V., L. A. Bastidas, S. Sorooshian, W. J. Shuttleworth, and Z. L. Yang, 1999: Parameter estimation of a land surface scheme using multicriteria methods. J. Geophys. Res., 104, 19 491–19 503, doi:10.1029/1999JD900154.
- —, H. Kling, K. K. Yilmaz, and G. F. Martinez, 2009: Decomposition of the mean squared error and NSE performance criteria: Implications for improving hydrological modelling. J. Hydrol., 377, 80–91, doi:10.1016/j. jhydrol.2009.08.003.
- Haughton, N., G. Abramowitz, A. Pitman, and S. J. Phipps, 2015: Weighting climate model ensembles for mean and variance estimates. *Climate Dyn.*, 45, 3169–3181, doi:10.1007/ s00382-015-2531-3.
- Henderson-Sellers, A., K. McGuffie, and A. J. Pitman, 1996: The Project for Intercomparison of Land-Surface Parametrization Schemes (PILPS): 1992 to 1995. *Climate Dyn.*, **12**, 849–859, doi:10.1007/s003820050147.

- Holtslag, A. A. M., G. J. Steeneveld, and B. J. H. Van de Wiel, 2007: Role of land-surface temperature feedback on model performance for the stable boundary layer. *Atmospheric Boundary Layers*, Springer, 205–220, doi:10.1007/978-0-387-74321-9_14.
- Ingwersen, J., and Coauthors, 2011: Comparison of Noah simulations with eddy covariance and soil water measurements at a winter wheat stand. *Agric. For. Meteor.*, **151**, 345–355, doi:10.1016/j.agrformet.2010.11.010.
- Kessomkiat, W., H.-J. H. Franssen, A. Graf, and H. Vereecken, 2013: Estimating random errors of eddy covariance data: An extended two-tower approach. *Agric. For. Meteor.*, **171–172**, 203–219, doi:10.1016/j.agrformet.2012.11.019.
- Koster, R. D., and Coauthors, 2004: Regions of strong coupling between soil moisture and precipitation. *Science*, **305**, 1138– 1140, doi:10.1126/science.1100217.
- —, and Coauthors, 2006: GLACE: The Global Land–Atmosphere Coupling Experiment. Part I: Overview. J. Hydrometeor., 7, 590–610, doi:10.1175/JHM510.1.
- —, Z. Guo, R. Yang, P. A. Dirmeyer, K. Mitchell, and M. J. Puma, 2009: On the nature of soil moisture in land surface models. *J. Climate*, **22**, 4322–4335, doi:10.1175/2009JCLI2832.1.
- Manabe, S., 1969: Climate and the ocean circulation: I. The atmospheric circulation and the hydrology of the earth's surface. *Mon. Wea. Rev.*, 97, 739–774, doi:10.1175/1520-0493(1969)097<0739: CATOC>2.3.CO;2.
- Masson, D., and R. Knutti, 2011: Climate model genealogy. *Geophys. Res. Lett.*, **38**, L08703, doi:10.1029/2011GL046864.
- Mauder, M., and T. Foken, 2006: Impact of post-field data processing on eddy covariance flux estimates and energy balance closure. *Meteor. Z.*, **15**, 597–609, doi:10.1127/ 0941-2948/2006/0167.
- Meehl, G. A., C. Covey, T. Delworth, M. Latif, B. McAvaney, J. Mitchell, R. Stouffer, and K. Taylor, 2007: The WCRP CMIP3 multi-model dataset: A new era in climate change research. *Bull. Amer. Meteor. Soc.*, 88, 1383–1394, doi:10.1175/ BAMS-88-9-1383.
- Mendoza, P. A., M. P. Clark, M. Barlage, B. Rajagopalan, L. Samaniego, G. Abramowitz, and H. Gupta, 2015: Are we unnecessarily constraining the agility of complex processbased models? *Water Resour. Res.*, **51**, 716–728, doi:10.1002/ 2014WR015820.
- Monteith, J. L., and M. H. Unsworth, 1990: *Principles of Envi*ronmental Physics. Butterworth-Heinemann, 291 pp.
- Pitman, A. J., 2003: The evolution of, and revolution in, land surface schemes designed for climate models. *Int. J. Climatol.*, 23, 479–510, doi:10.1002/joc.893.
- —, and Coauthors, 1999: Key results and implications from phase 1(c) of the Project for Intercomparison of Land-Surface Parametrization Schemes. *Climate Dyn.*, **15**, 673–684, doi:10.1007/s003820050309.
- Schymanski, S. J., M. L. Roderick, M. Sivapalan, L. B. Hutley, and J. Beringer, 2007: A test of the optimality approach to modelling canopy properties and CO₂ uptake by natural vegetation. *Plant Cell Environ.*, **30**, 1586–1598, doi:10.1111/ j.1365-3040.2007.01728.x.
- Seneviratne, S. I., and Coauthors, 2013: Impact of soil moisture– climate feedbacks on CMIP5 projections: First results from the GLACE-CMIP5 experiment. *Geophys. Res. Lett.*, 40, 5212–5217, doi:10.1002/grl.50956.
- Stoy, P. C., and Coauthors, 2013: A data-driven analysis of energy balance closure across FLUXNET research sites: The role of landscape scale heterogeneity. *Agric. For. Meteor.*, **171–172**, 137–152, doi:10.1016/j.agrformet.2012.11.004.

- Taylor, K. E., R. J. Stouffer, and G. A. Meehl, 2012: An overview of CMIP5 and the experiment design. *Bull. Amer. Meteor. Soc.*, 93, 485–498, doi:10.1175/BAMS-D-11-00094.1.
- van den Hurk, B. J. J., M. J. Best, P. A. Dirmeyer, A. J. Pitman, J. Polcher, and J. Santanello, 2011: Acceleration of land surface model development over a decade of GLASS. *Bull. Amer. Meteor. Soc.*, 92, 1593–1600, doi:10.1175/BAMS-D-11-00007.1.
- Wei, J., P. A. Dirmeyer, Z. Guo, L. Zhang, and V. Misra, 2010: How much do different land models matter for climate

simulation? Part I: Climatology and variability. J. Climate, 23, 3120–3134, doi:10.1175/2010JCLI3177.1.

- Wilson, K., and Coauthors, 2002: Energy balance closure at FLUXNET sites. Agric. For. Meteor., 113, 223–243, doi:10.1016/S0168-1923(02)00109-0.
- Zaehle, S., and Coauthors, 2014: Evaluation of 11 terrestrial carbon-nitrogen cycle models against observations from two temperate Free-Air CO₂ Enrichment studies. *New Phytol.*, 202, 803–822, doi:10.1111/nph.12697.