

The water balance representation in Urban-PLUMBER land surface models

Article

Published Version

Creative Commons: Attribution 4.0 (CC-BY)

Open Access

Jongen, H. J., Lipson, M., Teuling, A. J., Grimmond, S.
ORCID: <https://orcid.org/0000-0002-3166-9415>, Baik, J.-J.,
Best, M., Demuzere, M., Fortuniak, K., Huang, Y., De Kauwe,
M. G., Li, R., McNorton, J., Meili, N., Oleson, K., Park, S. B.,
Sun, T. ORCID: <https://orcid.org/0000-0002-2486-6146>,
Tsiringakis, A., Varentsov, M., Wang, C., Wang, Z. H. and
Steeneveld, G. J. (2024) The water balance representation in
Urban-PLUMBER land surface models. *Journal of Advances in
Modeling Earth Systems*, 16 (10). e2024MS004231. ISSN
1942-2466 doi: 10.1029/2024MS004231 Available at
<https://centaur.reading.ac.uk/117958/>

It is advisable to refer to the publisher's version if you intend to cite from the work. See [Guidance on citing](#).

To link to this article DOI: <http://dx.doi.org/10.1029/2024MS004231>

Publisher: American Geophysical Union

All outputs in CentAUR are protected by Intellectual Property Rights law, including copyright law. Copyright and IPR is retained by the creators or other copyright holders. Terms and conditions for use of this material are defined in the [End User Agreement](#).

www.reading.ac.uk/centaur

CentAUR

Central Archive at the University of Reading

Reading's research outputs online



RESEARCH ARTICLE

10.1029/2024MS004231

Key Points:

- We evaluate the water balance in 19 urban land surface models (ULSM) from the Urban-PLUMBER project
- ULSMs capture the timing of water fluxes more accurately than their magnitude
- The water balance appears unclosed in 43% of the model runs (19 models at 20 sites)

Supporting Information:

Supporting Information may be found in the online version of this article.

Correspondence to:

G.-J. Steeneveld,
gert-jan.steeneveld@wur.nl

Citation:

Jongen, H. J., Lipson, M., Teuling, A. J., Grimmond, S., Baik, J.-J., Best, M., et al. (2024). The water balance representation in urban-PLUMBER land surface models. *Journal of Advances in Modeling Earth Systems*, 16, e2024MS004231. <https://doi.org/10.1029/2024MS004231>

Received 23 JAN 2024




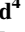








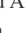
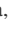


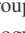
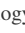

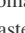
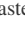
Accepted 26 AUG 2024

Author Contributions:

Conceptualization: H. J. Jongen, A. J. Teuling, G. J. Steeneveld
Data curation: H. J. Jongen, M. Lipson
Formal analysis: H. J. Jongen
Funding acquisition: H. J. Jongen, A. J. Teuling, G. J. Steeneveld
Investigation: M. Lipson, S. Grimmond, J.-J. Baik, M. Best, M. Demuzere, K. Fortuniak, Y. Huang, M. G. De Kauwe, R. Li, J. McNorton, N. Meili, K. Oleson, S.-B. Park, T. Sun, A. Tsiringakis, M. Varentsov, C. Wang, G. J. Steeneveld
Methodology: H. J. Jongen, M. Lipson, A. J. Teuling, S. Grimmond, G. J. Steeneveld
Project administration: H. J. Jongen, M. Lipson

© 2024 The Author(s). Journal of Advances in Modeling Earth Systems published by Wiley Periodicals LLC on behalf of American Geophysical Union. This is an open access article under the terms of the [Creative Commons Attribution License](#), which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

The Water Balance Representation in Urban-PLUMBER Land Surface Models

H. J. Jongen^{1,2} , M. Lipson³ , A. J. Teuling¹ , S. Grimmond⁴ , J.-J. Baik⁵ , M. Best⁶ , M. Demuzere^{7,8} , K. Fortuniak⁹ , Y. Huang¹⁰ , M. G. De Kauwe¹¹ , R. Li^{12,13} , J. McNorton¹⁴ , N. Meili^{15,16} , K. Oleson¹⁷ , S.-B. Park¹⁸ , T. Sun¹² , A. Tsiringakis^{2,19} , M. Varentsov²⁰ , C. Wang^{10,21} , Z.-H. Wang²² , and G. J. Steeneveld² 

¹Hydrology and Environmental Hydraulics, Wageningen University, Wageningen, The Netherlands, ²Meteorology and Air Quality, Wageningen University, Wageningen, The Netherlands, ³Bureau of Meteorology, Canberra, ACT, Australia, ⁴Department of Meteorology, University of Reading, Reading, UK, ⁵School of Earth and Environmental Sciences, Seoul National University, Seoul, South Korea, ⁶Met Office, Exeter, UK, ⁷Department of Geography, Urban Climatology Group, Ruhr-University Bochum, Bochum, Germany, ⁸B-Kode, Ghent, Belgium, ⁹Department of Meteorology and Climatology, Faculty of Geographical Sciences, University of Łódź, Łódź, Poland, ¹⁰School of Meteorology, University of Oklahoma, Norman, OK, USA, ¹¹School of Biological Sciences, University of Bristol, Bristol, UK, ¹²Institute for Risk and Disaster Reduction, University College London, London, UK, ¹³Department of Hydraulic Engineering, Tsinghua University, Beijing, China, ¹⁴European Centre for Medium-Range Weather Forecasts (ECMWF), Reading, UK, ¹⁵Department of Civil and Environmental Engineering, National University of Singapore, Singapore, Singapore, ¹⁶Future Cities Laboratory Global, Singapore-ETH Centre, Singapore, Singapore, ¹⁷U.S. National Science Foundation National Center for Atmospheric Research (NSF NCAR), Boulder, CO, USA, ¹⁸School of Environmental Engineering, University of Seoul, Seoul, South Korea, ¹⁹European Centre for Medium-Range Weather Forecasts (ECMWF), Bonn, Germany, ²⁰Faculty of Geography/Research Computing Center, Lomonosov Moscow State University, Moscow, Russia, ²¹Department of Geography and Environmental Sustainability, University of Oklahoma, Norman, OK, USA, ²²School of Sustainable Engineering and the Built Environment, Arizona State University, Tempe, AZ, USA

Abstract Urban Land Surface Models (ULSMs) simulate energy and water exchanges between the urban surface and atmosphere. However, earlier systematic ULSM comparison projects assessed the energy balance but ignored the water balance, which is coupled to the energy balance. Here, we analyze the water balance representation in 19 ULSMs participating in the Urban-PLUMBER project using results for 20 sites spread across a range of climates and urban form characteristics. As observations for most water fluxes are unavailable, we examine the water balance closure, flux timing, and magnitude with a score derived from seven indicators expecting better scoring models to capture the latent heat flux more accurately. We find that the water budget is only closed in 57% of the model-site combinations assuming closure when annual total incoming fluxes (precipitation and irrigation) fluxes are within 3% of the outgoing (all other) fluxes. Results show the timing is better captured than magnitude. No ULSM has passed all water balance indicators for any site. Models passing more indicators do not capture the latent heat flux more accurately refuting our hypothesis. While output reporting inconsistencies may have negatively affected model performance, our results indicate models could be improved by explicitly verifying water balance closure and revising runoff parameterizations. By expanding ULSM evaluation to the water balance and related to latent heat flux performance, we demonstrate the benefits of evaluating processes with direct feedback mechanisms to the processes of interest.

Plain Language Summary Urban environments have their own local climates including typically higher nocturnal temperatures compared with rural areas. Ideally, modeling cities should capture their influences on the atmosphere above them. As the energy and water balances are linked by evaporation, a good water balance representation will support a good energy balance simulation. Focusing on the water balance, we find the water balance in models could be improved by paying attention to closure and runoff.

1. Introduction

The impact of urbanization on the local climate and hydrology has sparked scientists' interest and inspired research for centuries (e.g., Howard, 1833; Oke, 1982; Fletcher et al., 2013; Hamdi et al., 2020). With the increasing population in cities (United Nations, 2018) more people are impacted by increased heat stress and flooding (Botzen et al., 2020; Gasparrini et al., 2017; Heaviside et al., 2016; Zhou et al., 2019). Spatial

Software: H. J. Jongen, M. Lipson
Supervision: A. J. Teuling,
G. J. Steeneveld
Visualization: H. J. Jongen
Writing – original draft: H. J. Jongen,
M. Lipson, A. J. Teuling, S. Grimmond, J.-
J. Baik, M. Best, M. Demuzere,
K. Fortuniak, Y. Huang, M. G. De Kauwe,
R. Li, J. McNorton, N. Meili, K. Oleson,
S.-B. Park, T. Sun, A. Tsiringakis,
M. Varentsov, C. Wang, G. J. Steeneveld

morphological heterogeneity and human interactions make understanding the urban climate challenging (Demuzere et al., 2022; Koopmans et al., 2020; Kotthaus & Grimmond, 2014a; Sun et al., 2018), but weather and climate models need to include the effects of urban areas, as they locally exacerbate extreme events (Hertwig et al., 2020; Oleson et al., 2008; Ronda et al., 2017). Examples are increased flooding due to high impervious fractions (Zhou et al., 2019) and increased heat stress during heat waves resulting from reduced evaporation (Lemonsu et al., 2015; Li et al., 2019). Therefore, models need to capture the impact of urban areas on their climate.

Researchers have developed, evaluated, and improved Urban Land Surface Models (ULSMs) simulating the interaction of the urban surface with the atmosphere. Coupled with a numerical weather prediction or climate model, ULSMs serve as a lower boundary condition and improve the model performance for urban environments (Tewari et al., 2007). ULSMs make different simplifying assumptions regarding urban geometry: a single homogeneous, impervious slab; multiple, individually homogeneous slabs; two-dimensional canyons; or 3D streets with individual buildings (Grimmond et al., 2009). These models also differ in whether and how they include physical processes like anthropogenic heat, irrigation, and snow processes (Lipson et al., 2024). To evaluate their performance, individual models are compared with observations (e.g., Grimmond & Oke, 2002; Hamdi & Schayes, 2007; Kravynhoff & Voogt, 2007; Porson et al., 2010; Ross & Oke, 1988). Although these individual evaluations were sometimes based on the same observations (Grimmond et al., 2009), the lack of a systematic approach prevented consistent comparison of the schemes. To compare the wide variety of models, two successive comparison projects applied a systematic approach. The first systematic comparison of ULSMs generally followed the PILPS protocol (Project for Intercomparison of Land surface Parameterization Schemes, Henderson-Sellers et al. (1996)), hence PILPS-Urban (Grimmond et al., 2010, 2011). Individual modelers received meteorological input and surface characteristics to enable them to run their models. In total, 32 models completed simulations for a site in Vancouver and one in Melbourne. Grimmond et al. (2011) concluded that increased model complexity did not necessarily benefit model performance.

The second intercomparison, Urban-PLUMBER (Lipson et al., 2024), assesses 30 models initially at the PILPS-Urban Melbourne site and adopts benchmarks following the PLUMBER project (Best et al., 2015). Benchmarks serve as a relative reference, to which models are compared to assess whether a cohort performs better (or not) than the benchmark and if input information is utilized effectively. Urban-PLUMBER is extended to the 20 sites presented by Lipson et al. (2022b) in the second phase (Lipson et al., 2023). The Urban-PLUMBER models outperform the PILPS-Urban ones for the sensible and latent heat flux. Some models representing two-dimensional canyons now perform nearly as well as one and two-tile models after efforts to improve hydrology and vegetation representation. However, models with complex urban geometry often still have relatively simple hydrology and vegetation and perform less well overall suggesting the representation of hydrology and vegetation requires more attention (Lipson et al., 2024).

Although PILPS-Urban and Urban-PLUMBER conclude vegetation and hydrology are important for model performance, neither project evaluates the water balance explicitly. The water balance satisfies the conservation of mass (Lavoisier, 1789) in the same way the energy balance satisfies the conservation of energy (Châlet, 1740). The conservation of energy is forced in many ULSMs to prevent the energetic state of the model from drifting and the consequential, long-term bias in the modeled surface fluxes (Grimmond et al., 2010). Closure is achieved by either updating the surface temperatures based on the residual energy or restricting the turbulent heat flows to the available energy (Grimmond et al., 2010). Both PILPS-Urban and Urban-PLUMBER test whether models close the energy balance, but have not verified the numerical closure of the water balance. Similar to the energy balance, an unclosed water balance can result in model biases and consequential drifting. These biases may in turn affect the energy balance, as the energy and water balance are linked through evapotranspiration (ET), the mass counterpart of the latent heat flux (Q_E). This direct link implies errors and/or biases in one balance will affect the model's skill for the other balance. Recently, Yu et al. (2022) showed the hydrology in a coupled ULSM has the potential to improve the Q_E , humidity, and air temperature with impacts up into the boundary layer (~ 1 km). ET/Q_E has been amongst the most challenging fluxes for ULSMs from the first assessment (Ross & Oke, 1988) until now (Lipson et al., 2024). Given the link to the energy balance, we hypothesize closing the water balance will improve model performance for the energy balance fluxes.

However, the water balance cannot be directly assessed because of a lack of observations at the appropriate spatiotemporal scales at this time. While precipitation is measured routinely in many urban locations with rain

Table 1

Overview of the Seven Indicators That are Linearly Combined in the UWBR Score, Which is Used to Evaluate the Urban Water Balance Representation in ULSMs

Water balance flux	Indicator	Description	Timescale	Criterion	Equation
All	I_A	Closure of the annual water balance assesses relative to the precipitation plus irrigation	Annual	<0.03	$\left \frac{P + I - (R + ET + \Delta S)}{P + I} \right $
ET	$I_{ET,m}$	Modeled cumulative ET normalized by the benchmark ET (ET_{bench}) over the whole model period	Modeled period	Within benchmark uncertainty*	$\frac{ET_{model}}{ET_{bench}}$
	$I_{ET,t}$	Similarity of ET recession timescale distribution between model and observations from the whole model run	Modeled period	$p < 0.05$	Kolmogorov-Smirnov test (Chakravarti et al., 1967)
ΔS	$I_{S,m}$	Range over the whole model run in stored water for both the modeled explicit and implicit water storage compared to water storage capacity	Modeled period	$< (50\% \text{ of soil volume} + 3 \text{ mm interception})$	$\Delta S_{model,max} - \Delta S_{model,min}$ (Equation 2) and $\Delta S_{max} - \Delta S_{min}$ (Equation 1)
	$I_{S,t}$	Coefficient of determination (R^2) between changes in explicit and implicit modeled water storage over the whole model period	Modeled period	> 0.9	$1 - \frac{\sum_{i=1}^n (\Delta S_i - \Delta \hat{S}_{model,i})^2}{\sum_{i=1}^n \Delta S_i - \Delta \hat{S}_{model,i}}$
R_s	$I_{R,m}$	Curve number (CN) from modeled runoff events and from site characteristics	Event	Within CN uncertainty*	$CN = \frac{1000}{S-10}$ (Section 2.1.4)
	$I_{R,t}$	Mean lag (hours) between centre of mass from precipitation and surface runoff of all events	Event	$< 1 \text{ hour}$	$R_{s,centroid} - P_{centroid}$

Note. The criterion indicates what needs to be achieved to assign a value of 1 to the indicator or 0.5 per test in the case of $I_{S,m}$. The uncertainty criteria (*) are discussed in Sections 2.1.2 and 2.1.4. The notation in the equations is defined in the corresponding subsections of Section 2.1. The details on all indicators can be found in Section 2.1.

gauges and rain radars, runoff, irrigation, and changes in water storage are not. Q_E (ET) observations from eddy-covariance systems have substantial gaps introduced in the quality control process (Feigenwinter et al., 2012) that rejects more data close to rain events (Grimmond, 2006). Runoff is occasionally measured in urban catchments (Berthier et al., 1999; Walsh et al., 2005), but a challenge is posed by the difference in the source area of observations for runoff and eddy-covariance techniques (Grimmond & Oke, 1986, 1991; Hellsten et al., 2015). External water use, often irrigation, further complicates the water balance in cities, as it mainly occurs at the micro-scale (e.g., garden irrigation). This scale can only be inferred from neighborhood piped water supply observations and water use surveys or estimated from weather, vegetation, and soil type (Grimmond & Oke, 1986; Kokkonen et al., 2018; Mitchell et al., 2001; Zeisl et al., 2018). Tree roots penetrate (sewer) pipes causing damage (Randrup et al., 2001) and simultaneously taking out water, which is an unobserved term. Lastly, measuring the water storage change is logistically difficult, as this requires the state of each individual element contributing to water storage in the city, such as soil moisture, interception, groundwater, and surface water. Thus, a direct comparison of a full set of water balance observations is extremely challenging and an alternative approach is needed.

Here, we develop an alternative approach to evaluate the representation and dynamics of the water balance in ULSMs. To examine the water balance closure, we propose an UWBR (urban water balance representation) score. The score combines seven indicators assessing: water balance closure (1 indicator), ET (2), water storage dynamics (2), and surface runoff (2). The UWBR score is applied, given a lack of observations, to rank models' capability to accurately capture different aspects of the water balance. Assessing the score of 19 Urban-PLUMBER ULSMs with a complete water balance representation helps to identify model improvement possibilities. The water balance representation is compared with the turbulent heat fluxes model skill since we expect a better water balance representation should improve simulated latent heat fluxes.

2. Methods

2.1. Urban Water Balance Representation (UWBR) Score

The UWBR score is a linear sum of seven indicators of a good water balance, which are assigned a value of one if a specified threshold is passed (Table 1), except the $I_{S,m}$ indicator, for which both sub-metrics are assigned 0.5 if passed. No weights are assigned, as these cannot be determined objectively. The UWBR score is compared with

the model performance for the latent heat flux assessed with metrics capturing different characteristics (Willmott, 1982) that are not entirely independent:

- Absolute mean bias error ($|MBE|$) assesses the bias providing insight into how well the quantities of the latent heat flux are modeled.
- Coefficient of determination (R^2) captures the consistency of the timing as R^2 decreases with a shift in a quasiperiodic signal like the latent heat flux.
- Normalized standard deviation (σ_{norm} , σ_{model} divided by $\sigma_{\text{observations}}$) compares the variability, which is dominated by the daily cycle in the case of the latent heat flux.
- Systematic Mean Absolute Error (MAE_s) indicates the average error. The systematic error is separated from the unsystematic error similar to the approach presented by Willmott (1982) for the root mean square error. This separation allows us to distinguish between systematic and random errors.
- Unsystematic Mean Absolute Error (MAE_u) assesses how well the erratic behavior is captured.

Before the individual indicators are introduced, we define two ways to calculate water storage from the model output based on either the water storage term (explicit) or the other terms of the water balance combined (implicit). Assuming that the net change in water stored in a “catchment” or a model grid (ΔS) can be derived from the difference between the incoming and outgoing water fluxes, then the implicit water storage is:

$$\Delta S = P + I - (R + ET) \quad (1)$$

where P is precipitation, I irrigation, and R runoff. R represents both the surface (R_s) and the subsurface (R_{sub}) runoff. When ΔS is calculated from the fluxes on the right-hand side of Equation 1, we refer to this as the implicit water storage. The second approach determines the net storage change (ΔS) based on the modeled storage components following the urban water balance (Grimmond & Oke, 1986). The storage components should account for the water storage above and below ground, such as the interception, water bodies, and groundwater. The components included depend on the model conceptualization. Here, we refer to the storage represented in the model as the explicit water storage (ΔS_{model}):

$$\Delta S_{\text{model}} = \Delta S_{\text{soil}} + \Delta S_{\text{intercept}} + \Delta S_{\text{snow}} \quad (2)$$

where ΔS_{soil} is storage change in the soil moisture, $\Delta S_{\text{intercept}}$ storage change in the interception storage, and ΔS_{snow} storage change in the snow cover. Depending on the model, ΔS_{soil} considers soil moisture below the impervious and pervious fraction. In the case a model does not consider soil moisture below the impervious fraction, ΔS_{soil} is adjusted accordingly. When we refer to annual timescales, the analysis is performed on all time intervals of a year in the time series, that is, a new annual period starts at every timestep, after which a full year is modeled (e.g., NL-Amsterdam: 2018-05-01 19:00–2019-05-01 19:00, 2018-05-01 20:00–2019-05-01 20:00, etc.). Within this year, no gaps in the model data allow all timesteps to be used. This method maximizes the use of available data and eliminates the influence of choosing a specific annual period like the calendar or hydrological year.

2.1.1. Water Balance Closure

Water balance closure assumes that all fluxes add up to zero for the time and space under consideration (here $\sim 1 \text{ km}^2$ and 1 year):

$$P + I - (R + ET + \Delta S_{\text{model}}) = 0 \quad (3)$$

where ΔS corresponds to the explicit water storage in the model (Equation 2) to prevent closure resulting from calculating the storage change based on the fluxes. Three models (8, 16, and 17) model groundwater interaction, which is not included in the model output. We examine the annual water balance closure with the annual total fluxes. Closure should also occur at every timestep, however, we were unable to undertake this more stringent check because interception storage was modeled but mostly unreported by modelers (all 19 models modeled, only 3 reported). Assuming an interception storage capacity of over 0.5 and up to 3 mm (Carlyle-Moses et al., 2020; Klaassen et al., 1998; Wouters et al., 2015), this storage can be filled in a single (half-)hourly timestep. At a single

(half-)hourly timestep, 0.5 mm is a non-negligible lack of closure but it is less critical at the annual scale. Equation 3 is normalized by annual precipitation plus irrigation to enable comparison between sites with a range of precipitation regimes.

The water balance closure indicator (I_A , Table 1) assesses if the total sum of all fluxes (including storage) is less than 3% from $P + I$. The 3% threshold allows for non-closure due to interception storage data not being provided in the model output, errors arising in latent heat flux unit conversion, or numerical model errors. According to the literature, interception storage amounts to 0.5–3 mm explaining a non-closure of up to 0.5% when it is not provided (Carlyle-Moses et al., 2020; Klaassen et al., 1998; Wouters et al., 2015). Converting the latent heat flux to ET can result in variations up to 2% depending on temperature and snow effects (Bringfelt, 1986; Petrucci et al., 2010). Not all models correct for these effects. To account for numerical model errors arising from discretization and time stepping (MacKay et al., 2022), we allow deviations of up to 0.5%.

2.1.2. Evapotranspiration (ET)

The two ET indicators address the magnitude and timing. The non-randomly distributed gaps in ET observations prevent direct comparison of total modeled ET (ET_{model}) over a model period. Thus, we use one of the Lipson et al. (2024) benchmark models. This allows a total ET to be obtained without gaps. The Lipson et al. (2024) benchmark model (ET_{bench}) is derived using multivariate ordinary least squares regressions with a K-means clustering approach. The K-means clustering approach is trained in-sample using 81 clusters on four variables: incoming shortwave radiation, air temperature, relative humidity, and wind speed (KM4-IS-SWdown-Tair-RH-Wind in Lipson et al., 2024). To reduce the hourly MBE, wind speed is omitted at both Helsinki sites. At all sites, the MBE is below 1 W m^{-2} and at most sites below 0.1 W m^{-2} evaluated against available data.

Therefore, ET_{bench} is assumed to provide a reasonable estimate of the total ET flux over the model run for the $I_{ET,m}$ indicator (Table 1). We compare in Q_E units rather than ET , eliminating unit conversions and calculate the cumulative ET flux uncertainty from the benchmark based on (a) the benchmark MBE multiplied by the run duration, and (b) lack of energy balance closure associated with eddy-covariance observations (Foken et al., 2012; Franssen et al., 2010; Mauder et al., 2020). The lack of energy closure is calculated by the net all-wave radiation minus the sum of the turbulent heat fluxes. The storage and anthropogenic heat fluxes are not observed, which prevents constraining the turbulent heat fluxes with energy balance closure. If a lack of closure occurs, the unexplained energy over the whole model run is split between Q_E and the sensible heat flux (Q_H) according to the Bowen ratio based on the benchmark fluxes (Hirschi et al., 2017; Mauder et al., 2020; Twine et al., 2000):

$$Q_{E,uncertainty} = \frac{1}{1+B}(R_{net} - Q_E - Q_H) \text{ with } B = Q_H/Q_E \quad (4)$$

where B is the Bowen ratio and R_{net} the net radiation. To this $Q_{E,uncertainty}$, the benchmark uncertainty is added. The benchmark uncertainty is the MBE of the benchmark multiplied by the run duration. A model run passes $I_{ET,m}$ when ET_{model} falls within the uncertainty of ET_{bench} .

The timing of modeled ET is assessed assuming exponential ET recession after rainfall based on the recession timescale estimated following the Jongen et al. (2022) methodology. This methodology considers only the first 10 days to exclude the influence of longer dry periods and irrigation. A daily timescale analysis circumvents observational gaps. Model and observations are assessed if they have the same distribution for the recession timescale with a Kolmogorov-Smirnov test (Chakravarti et al., 1967). The $I_{ET,t}$ indicator is assigned a value of 1 when the p-value is below 0.05.

2.1.3. Water Storage

Indicator $I_{S,m}$ evaluates the water storage by comparing the modeled explicit and implicit water storage ranges (Section 2.1) over the analysis period with respect to the estimated water storage capacity. According to the literature, soil water storage capacity is maximally half the soil depth for all soil types (Saxton et al., 1986). The maximum is set as a storage capacity that models should not exceed rather than a realistic value. As urban soils are frequently disturbed making them spatially heterogeneous, reliable maps are rarely available (Van de Vijver et al., 2020). As the modeled soil depth depends on the model run, the soil water storage capacity is calculated for

each separately. To account for interception storage, 3 mm is added to the estimated water storage capacity based on tree and impervious interception observations (Carlyle-Moses et al., 2020; Klaassen et al., 1998; Wouters et al., 2015). The two models not including soil moisture do not pass the first check of this indicator and are only evaluated based on the implicit water storage (Table 2). Other models receive 0.5 score when either the modeled explicit or implicit water storage range falls within the estimated water storage capacity (or 1 for both).

Indicator $I_{S,t}$ quantifies the internal temporal consistency between the change in explicit (Equation 2) and the implicit (Equation 1), which should be indicating the same flux. The coefficient of determination R^2 (Willmott, 1982) is calculated using storage changes using 30-min (or 60-min) model output depending on the site forcing data. This metric equals one if the timing between two fluxes is similar ($R^2 > 0.9$) independent of the flux bias, unlike other indicators (e.g., I_A). The two models without soil moisture output are assigned a value of 0 for $I_{S,t}$ as their performance could not be evaluated.

2.1.4. Surface Runoff (R_s)

Indicator $I_{R,m}$ assesses the R_s magnitude relating total event precipitation to R_s (Figure 1a). Without runoff observations, curve numbers (CN) are derived to evaluate modeled total event R_s (Cronshey et al., 1985) based on the relation between the total event precipitation (P_e) and the total event R_s (R_e):

$$R_e = \frac{(P_e - 0.2S)^2}{P_e + 0.8S} \text{ with } S = \frac{1000}{CN} - 10 \quad (5)$$

where S is the potential maximum retention. To determine when precipitation events are independent, the autocorrelation of precipitation events is examined. A dry period of 5 hours (Figure S1 in Supporting Information S1) is assumed across all sites, which is consistent with Wenzel Jr and Voorhees (1981). This dry period is several hours longer than the expected runoff response time preventing events from influencing each other (Berne et al., 2004; Morin et al., 2001; Yao et al., 2016). To exclude snow events, the analysis includes only events with a minimum air temperature above 0°C. For each model run, ordinary least squares is used with the R_e and P_e data to estimate S (Figure 1b) from which the CN is derived (Equation 5). During this process, the variance and standard deviation of S are calculated from the variance in the data points. The standard deviation of CN follows from this and is used as the uncertainty estimate from the models.

For each site, the CN is estimated using a linear interpolation of a look-up table considering the impervious fraction within the eddy-covariance footprint (Cronshey et al., 1985). Given soil texture influences CN , sand fraction (Brakensiek & Rawls, 1983; Nachtergaele, 2001) obtained from a global data set (OpenLandMap, (Hengl, 2018)) is used to constrain CN . Given the uncertainty of urban soil maps, using sand fraction is a repeatable way to assign the most uncertainty to the CN look-up tables, assuming a one-third change of CN from a one-level change in soil texture in either direction. If the site CN , including its uncertainty, overlaps with the model CN including its uncertainty, $I_{R,m}$ is assigned a value of 1.

Indicator $I_{R,t}$ addresses the rainfall- R_s response times (Leopold, 1968). The lag time is calculated as the difference between centroids of rainfall ($P_{centroid}$) and R_s ($R_{centroid}$) for the same events as the CN calculations (Figure 1a). Long-tail rainfall events are excluded when the $R_{centroid}$ comes before the $P_{centroid}$. As eddy-covariance systems have a footprint on the sub-square-kilometer scale (Feigenwinter et al., 2012), lag time is expected to be much faster than 30–60 min (Berne et al., 2004; Morin et al., 2001; Yao et al., 2016), which is the model output resolution (Lipson et al., 2024). Therefore, the mean lag time needs to be less than 1 hour. The mean is preferred over the median to also pinpoint models that occasionally have long lag times that would not affect the median. Lag times of intermittent precipitation-runoff events will only decrease, as storages are already (partly) filled by earlier precipitation. Dry periods of less than 5 hours should also have lag times of less than 1 hour.

2.2. Models

The present study anonymously analyzes the water balance outputs from 19 Urban-PLUMBER ULSMs (Table 2). Other Urban-PLUMBER ULSMs did not submit the necessary outputs to allow for a water balance assessment. The outputs are for 20 sites covering a range of climates, impervious fractions, and observational periods (Table 3). As two models did not run all sites, 377 runs are analyzed.

Table 2

Overview of the 19 Urban Land Surface Models in the Water Balance Analysis Based on Lipson et al. (2024)

Model	Urban geometry	Vegetation	Soil hydrology	Snow accumulation	Irrigation	Water balance closure check	Reference
ASLUMv2.0	Canyon	Grass	Multi-layer	No	No ^b	No ^c	Z.-H. Wang et al. (2013) C. Wang et al. (2021)
ASLUMv3.1	Canyon	Grass + trees	Multi-layer	No	No ^b	No ^c	Z.-H. Wang et al. (2013) C. Wang et al. (2021)
CABLE	Non-urban	Separate tiles	Multi-layer	Veg.	No	Yes	Kowalczyk et al. (2006) Y. P. Wang et al. (2011)
ECLand	Non-urban	Separate tiles	Multi-layer	Veg.	No	No ^d	Boussetta et al. (2021)
ECLand-U	Two-tile	Separate tiles	Multi-layer	Veg. + urban	No	No ^d	McNorton et al. (2021) Boussetta et al. (2021)
CLMU5	Canyon	Grass + shrubs	Multi-layer	Urban	No	Yes	Oleson and Feddema (2020)
JULES 1T	One-tile	Separate tiles	Multi-layer	Veg. + urban	No	Yes	Best et al. (2011)
JULES 2T	Two-tile	Separate tiles	Multi-layer	Veg. + urban	No	Yes	Best et al. (2011)
JULES MOR	Two-tile	Separate tiles	Multi-layer	Veg. + urban	No	Yes	Best et al. (2011)
Lodz-SUEB	One-tile	Lumped with urban	Multi-layer ^a	Veg. + urban	No	No	Fortuniak (2003)
Manabe 1T	One-tile	Manabe bucket	One-layer	Veg. + urban	No	No	Best et al. (2011) Manabe (1969)
Manabe 2T	Two-tile	Manabe bucket	One-layer	Veg. + urban	No	No	Best et al. (2011) Manabe (1969)
NOAH-SLAB	One-tile	Separate tiles	Multi-layer	Veg. + urban	No	No	Kusaka et al. (2001) Ek et al. (2003)
NOAH-SLUCM	Canyon	Separate tiles	Multi-layer	Veg. + urban	No	No	Kusaka et al. (2001) Ek et al. (2003)
SNUUCM	Canyon	Separate tiles	Multi-layer ^a	Veg.	No	No	Ryu et al. (2011) Ek et al. (2003)
SUEWS	Two-tile	Separate tiles	One-layer	Veg. + urban	No ^b	Yes	Järvi et al. (2011) Ward et al. (2016)
TERRA 4.11	One-tile	Separate tiles	Multi-layer	Veg.	No	No	Wouters et al. (2015) Schulz and Vogel (2020)
UCLEM	Canyon	Grass + shrubs	One-layer	Veg. + urban	Yes	No	Thatcher and Hurley (2012) Lipson et al. (2018)
UT&C	Canyon	Grass + shrubs + trees	Multi-layer	No	Yes	Yes	Meili et al. (2020)

^aTwo models did not provide soil moisture output. ^bThree models capable of simulating irrigation did not include it in their Urban-PLUMBER runs. ^cTwo models have an internal water balance closure check which does not hold in conditions with snowfall. ^dTwo models have an internal water balance closure check that was not used for the Urban-PLUMBER runs.

For each site, modelers were provided with the site characteristics and meteorological forcing with 10-year spin-up data (Lipson et al., 2022b). The spin-up period required to reach equilibrium varies per model, with some requiring many years to come to hydrological equilibrium with the forcing meteorology (Best & Grimmer, 2016; Yang et al., 1995). The 10 years of spin-up before the evaluation observations allowed the soil moisture stores to equilibrate with local conditions prior to analysis. ERA5 reanalysis data (Hersbach et al., 2020) are used to derive hourly forcing with bias-correction including diurnal and seasonal effects for each site (Lipson et al., 2022b).

Depending on site data, evaluation is undertaken with 30- or 60-min fluxes for periods varying between 148 and 1,827 days (average 912 days, Table 3). Similar to the Urban-PLUMBER protocol, to minimize human errors,

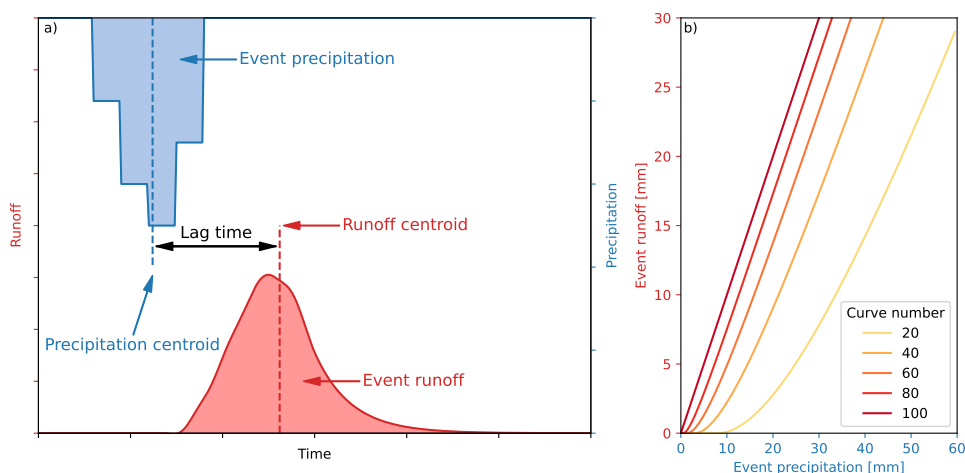


Figure 1. Illustration of surface runoff indicators ($I_{R,m}$ and $I_{R,t}$) showing (a) lag time between the precipitation centroid and surface runoff centroid for an illustrative event, and (b) CN values (Equation 5) derived from total event precipitation and surface runoff.

modelers received a preliminary analysis of the water balance to help identify major issues and were encouraged to update their results. This eliminated unit errors, added missing variables, and removed inactive soil moisture layers.

For this study, we harmonize the hydrological model output. If a model only provided Q_E (unit: $[W m^{-2}]$), it is converted to ET (unit: $[mm d^{-1}]$) using latent heat of vapourization accounting for air temperature (Brinck, 1986). When snow is present the latent heat of fusion is added to the latent heat of vapourization to acquire the latent heat of sublimation (Petrucci et al., 2010). In the forcing, precipitation is split into snowfall and rainfall. At only 30% of the sites, snowfall amounts to more than 10% of the precipitation. It is added as rainfall for one model without snow hydrology, while the two others do not account for this input. Irrigation is simulated in two models. For all other models, irrigation is assumed to be zero.

3. Results

The 19 ULSMs show a wide spread in the average yearly water fluxes at all 20 sites based on all 377 model runs (Figure 2). Overall, the model spread (whiskers, Figure 2) is often wider than the modeled ensemble mean flux (bars, Figure 2). Models show more variation in ET than in runoff. Sites with higher annual water input have more variability in model output fluxes, for example, the relatively high fluxes in KR-Jungnang and SG-TelokKurai compared to the lower yearly fluxes in PL-Lipowa and US-WestPhoenix.

3.1. Water Balance Closure

Although the annual mean model ensemble almost closes the water balance at most sites (Figure 2), most individual models do not close the water balance (Figure 3). Here, closure is assumed when the sum of all fluxes (Equation 3) is less than 3% of $P + I$. This occurs in 57% of the model runs (I_A , Figure 4). In 25% of the model runs, non-closure exceeds 10% of $P + I$. Closure is model-related as the bias is similar across sites for each model (Figure 3). Five models close the water balance in all runs, whereas four models account for 48% of unclosed model runs. Three models pass their internal water balance closure check but do not always pass this closure check possibly due to unreported, modeled water fluxes or inconsistencies in the way fluxes were reported. To assess the impact of model run length, the analysis is repeated with sites with more than 2 years of observations yielding similar results.

3.2. Evapotranspiration (ET)

Comparison of the modeled mean diurnal cycle of the ET (Figure 5) shows the highest inter-model spread at the peak of the diurnal cycle, with a range of 10%–600% of the model ensemble-mean flux. Along three sites with

Table 3

Model (Table 2) Outputs Are Analyzed for 20 Sites (Lipson et al., 2022b)

Country	City (site)	Name	Lat. (°)	Lon. (°)	Observed period (days)	Köppen-Geiger climate	LCZ	F_{imp}	z_d (m)	z_s (m)	Reference
Australia	Melbourne (Preston)	AU-Preston	−37.73	145.01	475	Cfb	6	0.62	8	40	Coutts et al. (2007a) Coutts et al. (2007b)
Australia	Melbourne (Surrey Hills)	AU-SurreyHills	−37.83	145.10	148	Cfb	6	0.54	8	38	Coutts et al. (2007a) Coutts et al. (2007b)
Canada	Vancouver (Sunset)	CA-Sunset	49.23	−123.08	1,827	Csb	6	0.68	3	25	Christen et al. (2011) Crawford and Christen (2015)
Finland	Helsinki (Kumpula)	FI-Kumpula	60.20	24.96	1,096	Dfb	Mix	0.46	6	31	Karsisto et al. (2016)
Finland	Helsinki (Torni)	FI-Torni	60.17	24.94	1,096	Dfb	2	0.77	15	60	Nordbo et al. (2013) Järvi et al. (2018)
France	Toulouse (Capitole)	FR-Capitole	43.60	1.45	375	Cfa	2	0.90	11	48	Masson et al. (2008) Goret et al. (2019)
Greece	Heraklion	GR-HECKOR	35.34	25.13	367	Csa	3	0.92	17	27	Stagakis et al. (2019)
Japan	Tokyo (Yoyogi)	JP-Yoyogi	35.66	139.68	1,461	Cfa	2	0.92	28	52	Hirano et al. (2015) Ishidoya et al. (2020)
South Korea	Seoul (Jungnang)	KR-Jungnang	37.59	127.08	825	Dwa	3	0.97	15	42	J.-W. Hong et al. (2020) S.-O. Hong et al. (2023)
South Korea	Cheongju (Ochang)	KR-Ochang	36.72	127.43	780	Dwa	5	0.47	4	19	J.-W. Hong et al. (2019) J.-W. Hong et al. (2020)
Mexico	Mexico City (Escandon)	MX-Escandon	19.40	−99.18	470	Cwb	2	0.94	8	37	Velasco et al. (2011) Velasco et al. (2014)
Netherlands	Amsterdam	NL-Amsterdam	52.37	4.89	652	Cfb	2	0.68	10	40	Steenefeld et al. (2020)
Poland	Łódź (Lipowa)	PL-Lipowa	51.76	19.45	1,827	Dfb	2	0.76	7	37	Pawlak et al. (2011) Fortuniak et al. (2013)
Poland	Łódź (Narutowicza)	PL-Narutowicza	51.77	19.48	1,827	Dfb	2	0.65	11	42	Fortuniak et al. (2006) Fortuniak et al. (2013)
Singapore	Singapore (Telok Kurau)	SG-TelokKurau	1.31	103.91	366	Af	3	0.85	7	24	Roth et al. (2017)
UK	London (King's college)	UK-KingsCollege	51.51	−0.12	638	Cfb	2	0.79	15	50	Kotthaus and Grimmond (2014a) Kotthaus and Grimmond (2014b) Bjorkegren et al. (2015)
UK	Swindon	UK-Swindon	51.58	−1.80	715	Cfb	6	0.49	4	13	Ward et al. (2013)
USA	Baltimore (Cub hill)	US-Baltimore	39.41	−76.52	1,826	Cfa	6	0.31	4	37	Crawford et al. (2011)
USA	Minneapolis	US-Minneapolis1	45.00	−93.19	1,093	Dfa	6	0.21	3	40	Peters et al. (2011) Menzer and McFadden (2017)
USA	Phoenix (West)	US-WestPhoenix	33.48	−112.14	382	Bwh	6	0.48	3	22	Chow et al. (2014) Chow (2017)

Note. Only wind directions coming from urban areas are included for the Minneapolis site. Characteristics include the local climate zone (LCZ, Stewart and Oke (2012), where two is compact mid-rise, 3 compact low-rise, five open mid-rise, and six open low-rise), impervious surface fraction (F_{imp}), displacement height (z_d), and eddy-covariance sensor height above ground level (z_s).

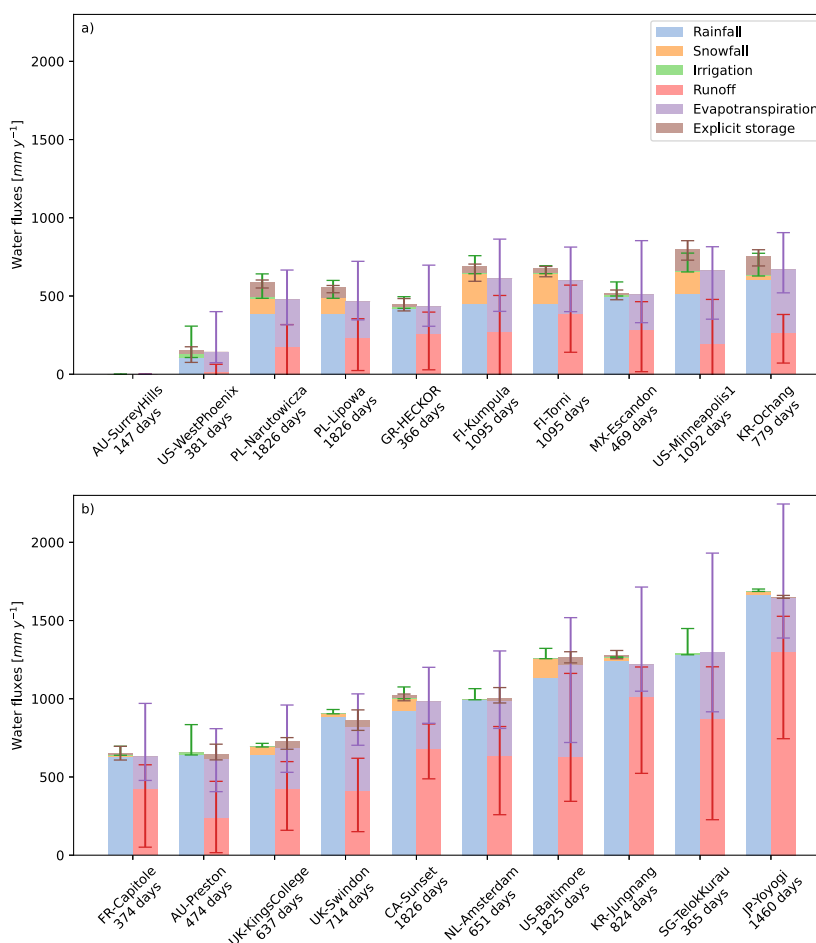


Figure 2. Ensemble mean (bars) and full range (minimum to maximum, whiskers) of the modeled annual water fluxes for all 20 sites ordered by increasing average annual precipitation. Explicitly modeled storage flux (Equation 2, brown) appears on the left if a net input and right if a net loss. Values are means of all complete years in a data set (e.g., NL-Amsterdam: 2018-05-01 19:00–2019-05-01 19:00, 2018-05-01 20:00–2019-05-01 20:00, etc.). AU-SurreyHills has less than a year of observations.

contrasting precipitation regimes (US-WestPhoenix, AU-Preston, and SG-TelokKurau), ET increases as expected at wetter sites. At US-WestPhoenix, all models but one underestimate peak ET . This underestimation likely results from the absence of irrigation in nearly all models, while irrigation is common at US-WestPhoenix (Templeton et al., 2018). The one overestimating model does not include irrigation. At the other two sites, around half the models underestimate ET (Figure 5). Although for these sites the model medians are better, the difficulty of capturing the correct flux magnitude is evident, as $I_{ET,m}$ is passed by only 26% of the model runs (Figure 4). No model passes this indicator at more than half of the sites.

After different rainfall events, daily ET decreases with varying timescales in both the observations and the models (Figure 6). The variation is higher amongst the modeled than the observed drydown. In contrast with the ET magnitude, the recession timescale shows no link with annual precipitation. $I_{ET,t}$ shows the ET recession timescale is captured correctly in 87% of the cases (Figure 4).

3.3. Water Storage

Not all models have explicit water storage values (Equation 2) that are equal to the implicit values (Equation 1, Figure 7), which is seen across all sites (not shown). However, the explicit water storage should reflect the implicit storage, as the explicit storage change is equal to the net of all water fluxes. For five models, the explicit storage

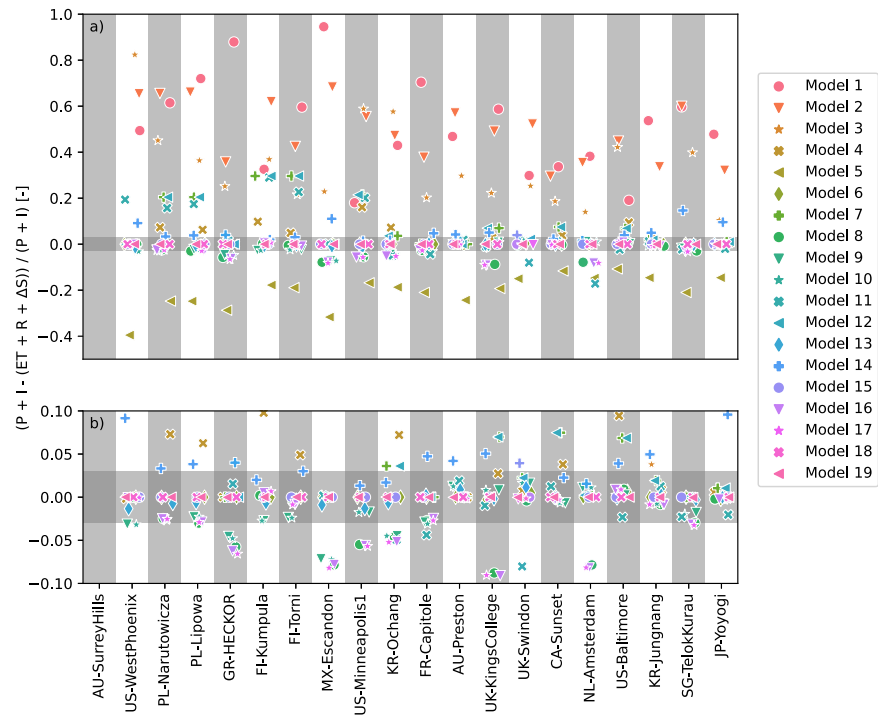


Figure 3. Annual water balance closure (Equation 3) per model (marker) at 20 sites (by increasing average annual precipitation). Models with indicator $I_A = 1$ (Table 1, horizontal shading) are shown in more detail in the lower panel (b).

change is equal to the implicit storage change at all sites. Minor differences occur in six models and large differences in six others. Two models have no differences at sites without snowfall (e.g., AU-Preston) but large differences at sites with snowfall (e.g., CA-Sunset). As these models do not account for the snowfall in the input

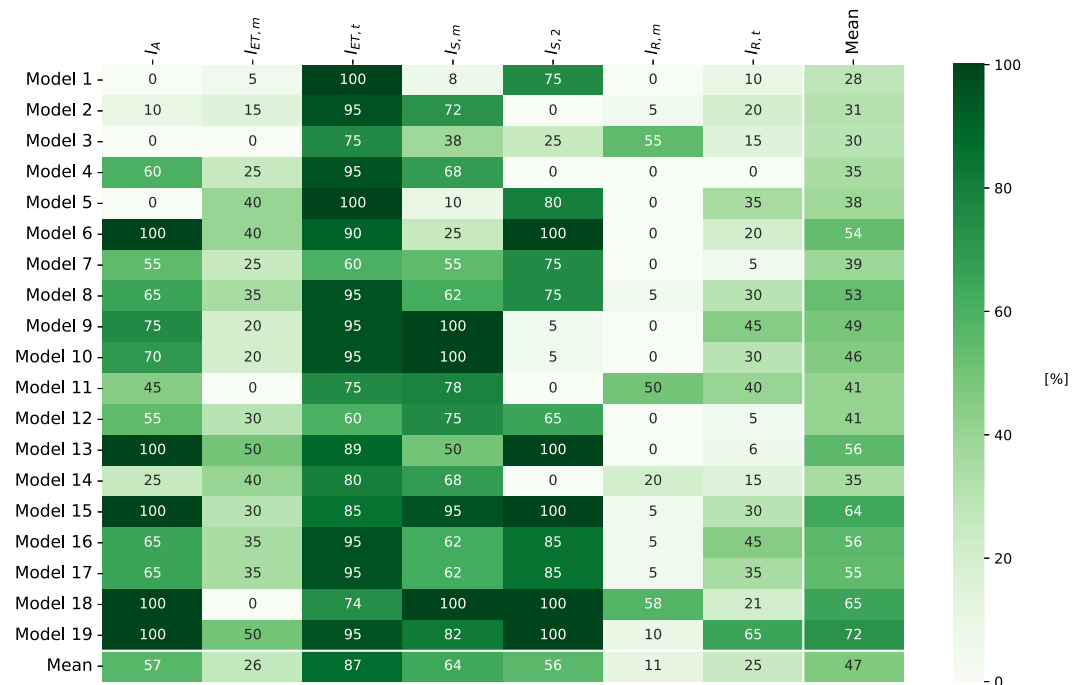


Figure 4. Overview of the indicators of the urban water balance representation (UWBR) score and constituent indicators (Table 1) over all sites. Means are corrected for missing model runs.

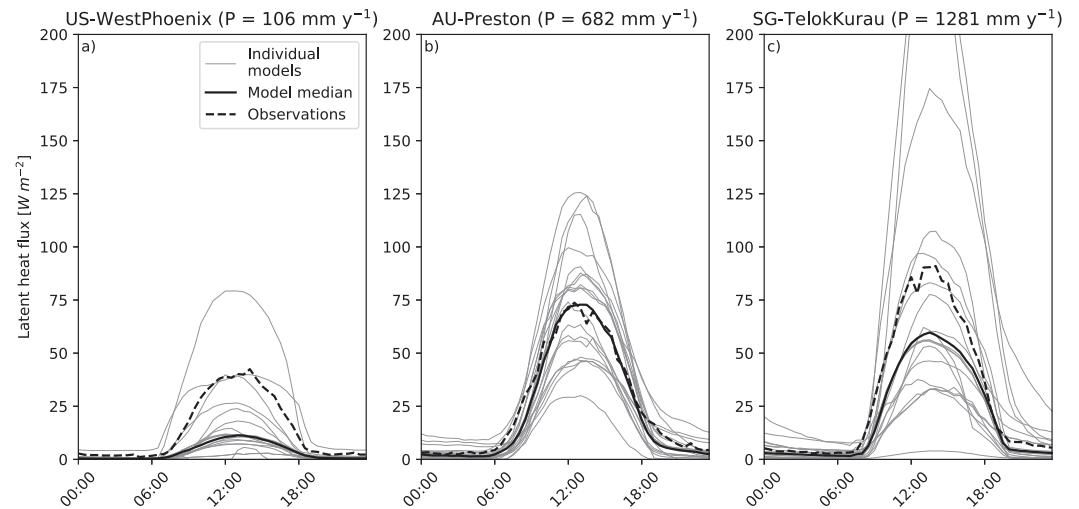


Figure 5. Illustration of modeled and observed (dashed) mean diurnal cycle of ET at three sites with contrasting annual precipitation increasing from left to right: (a) US-WestPhoenix, (b) AU-Preston, and (c) SG-TelokKurai. Note that the observations are direct latent heat flux observations from eddy-covariance systems and do not refer to ET_{bench} .

we see an increasing difference between the explicit and implicit water storage. The models with larger differences follow a seasonal cycle likely caused by non-restricted implicit water storage combined with restricted explicit water storage by soil storage capacity.

The range of modeled water storage exceeds the estimated site water storage capacity ($I_{S,m}$) in 64% of cases (Figure 4). Models 1 and 5 have the lowest score for this indicator, because they have an inconsistency between the inputs and outputs (Equation 3) causing non-closure of the water balance at nearly all sites. Three models never exceed the estimated water storage capacity.

How explicit relates to implicit water storage is linked to the individual models given the consistent results across sites (Figure 8). With magnitude represented by water balance closure, we focus on the timing by assessing the explicit relative to the implicit water storage (Figures 9a–9c). Model runs can have comparable directions but different patterns, for example, model 11 (Figure 9a), comparable patterns but different magnitudes of change, for example, model 9 (Figure 9b), or virtually no differences (e.g., model 18, Figure 9c). The explicit and implicit water storage changes (Figures 9d–9f) emphasize the difference in timing, which is why the indicator uses the R^2 of these derivatives. Only five models have virtually no differences and thus an R^2 of 1 (Figure 4). Over half of the models have R^2 greater than 0.9 indicating timing consistency ($I_{S,t}$, Figure 4).

3.4. Surface Runoff (R_s)

All models have surface runoff triggered by precipitation, but the precipitation event size causing R_s events differs between models (Figure 10). The model rather than the site seems to explain triggering event size despite the variation amongst sites in impervious fractions and precipitation regimes. This suggests that surface runoff parameterization may be critical. Thus, we find a large inter-model spread in the cumulative modeled R_s (Figure 2). One model is excluded as it does not output R_s separately from R_{sub} . Ten models show the expected increase of cumulative R_s with increasing site impervious fraction ($p > 0.05$, Wald test (Wald, 1943)), whereas nine models do not (Figure S2 in Supporting Information S1).

Only in 43 of the 337 model runs, the CN (curve number: Section 2.1.4) is captured correctly, passing $I_{R,m}$ (Figure 4), so all other model runs have no overlap with the site estimates (see Section 2.1.4). Three models capture the CN correctly for at least half of their model runs and are responsible for 32 of the successful model runs. Most models do not match event precipitation and R_s relation. Most models underestimate the CN relative to the site estimate (Figure S3 in Supporting Information S1). Underestimating the CN indicates a model is overestimating surface interception and/or soil infiltration, reducing R_s (Equation 5).

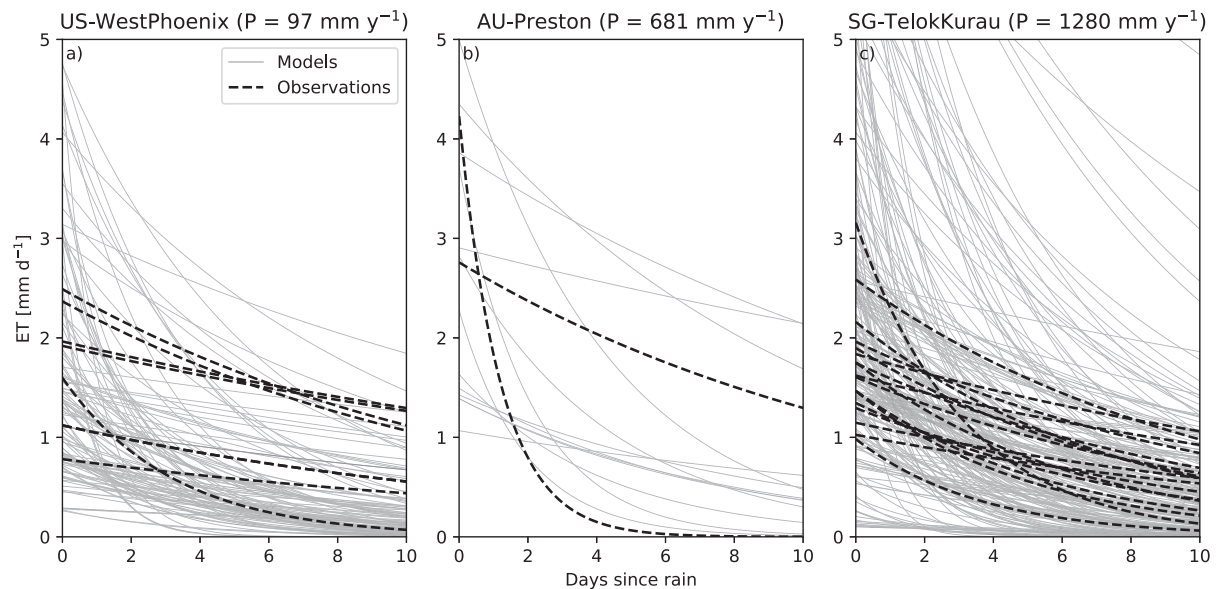


Figure 6. As Figure 5, but modeled (gray) and observed (black) daily ET following separate, individual rainfall events. Drydown events are selected based on their duration and data availability (see Jongen et al., 2022). Note that the observations are direct latent heat flux observations from eddy-covariance systems and do not refer to ET_{bench} .

One in four model runs accurately captures the fast R_s response in the lag time (Figure 4) with $I_{R,t}$ passed by 25% of the model runs. With very short lag times expected, only overestimates are simulated. Most lag times averaged per model run are less than 5 hours, but exceptionally they are over 100 hr. Average lag times per model run are shown in Figure S4 of Supporting Information S1.

3.5. Urban Water Balance Representation (UWBR) Score

Across all model runs, the mean UWBR score amounts to 3.3 out of the possible 7 (Figure 4). Although the overall pass rate across all indicators and models is 47%, pass rates strongly vary per indicator. Notably, 87% passes $I_{ET,t}$, while only 11% passes $I_{R,m}$. Pass rates also differ among models from 28% to 72%. Only one model run passes all indicators, while 10 model runs have a score of 6 out of 7. Model 19 accounts for five of these eleven high-scoring runs. If a model closes the water balance (I_A), it generally scores better on both storage indicators. In contrast, models with a high passing percentage for one ET indicator do not systematically score better for the other ET indicator. Overall, the ET timing ($I_{ET,t}$) is captured better than its cumulative magnitude ($I_{ET,m}$). A similar pattern is seen in the R_s indicators with the timing ($I_{R,t}$) captured slightly better than magnitude ($I_{R,m}$).

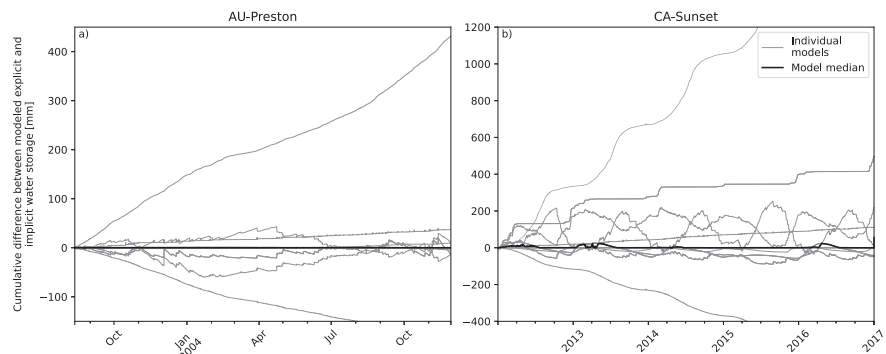


Figure 7. Cumulative difference between the explicit (Equation 2) and implicit water storage (Equation 1) at two representative sites for the entire model period for all models. Snowfall occurs at CA-Sunset, but not at AU-Preston. Some models are not visible as they are close to zero.

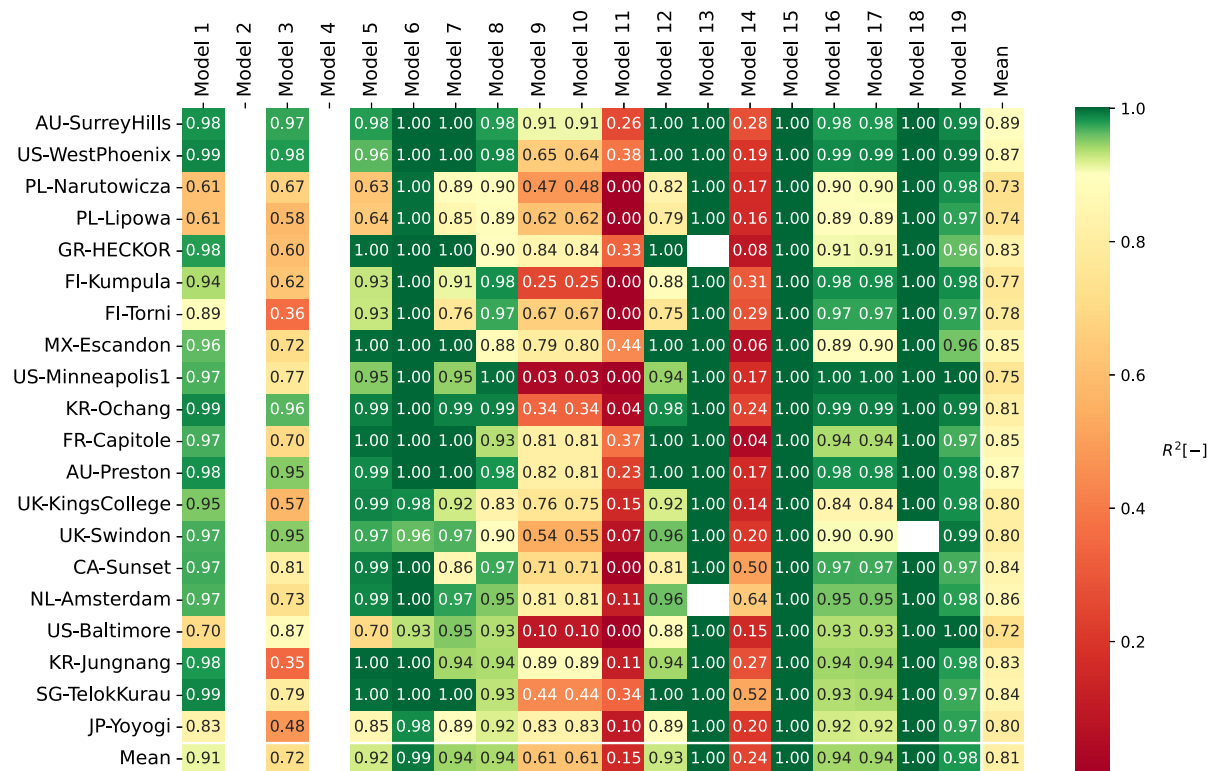


Figure 8. Coefficient of determination (R^2) between (half-)hourly explicit (Equation 2) and implicit water storage change (Equation 1) by model and site. Green indicates the 0.9 $I_{S,d}$ threshold (Table 1). Missing results are shown as white (i.e., cannot calculate explicit or implicit water storage change). Figure 9 may aid interpretation of R^2 values.

Generally, pass rates per indicator show a dependence on the model (Figure 4). This dependence is not found for sites (Figure S5 in Supporting Information S1). There is no relation evident between UWBR score and model approach (e.g., built surface, soil hydrology, Table 2), but the model is more influential than the site on UWBR score. As the Lipson et al. (2024) classification (Table 2) was not developed with the water balance representation as its original goal, further work would be needed to identify what model attributes are key to better UWBR score.

3.6. Linking the Water and Energy Balance

Surprisingly, models do not appear to capture any aspect of the latent heat flux more accurately if their UWBR score is higher. The UWBR score does not significantly correlate with better ranking on any of the four metrics evaluating the (half-)hourly modeled Q_E : the R^2 , σ_{norm} , MAE_s , and MAE_u ($p > 0.05$, Wald test, Figure S6 in Supporting Information S1). These correlations remain absent if one of the indicators is omitted from the analysis. The lack of correlation may be the result of the low number (11) of runs with a UWBR score higher than 5 (Figure 4) effectively reducing the UWBR score range. Given the lack of relations between the UWBR score and Q_E metrics, the Q_E is not better captured in model runs that pass more indicators of a realistic water balance representation, thus refuting our hypothesis that the urban water balance skill positively impacts simulated energy fluxes.

4. Discussion and Conclusions

This study assesses the water balance representation in 19 ULSMs from the Urban-PLUMBER project. It appears the water balance is not closed (within 3%) in 57% of the model-site runs. The considerable spread in water fluxes is as wide as the absolute flux magnitude at all sites. For both ET and R_s , the timing is captured better than the flux magnitude. Modeled explicit water storage dynamics (Equation 2) are inconsistent with the implicit water storage (Equation 1) in 44% of the models. Refuting our hypothesis, a better water balance representation does not result

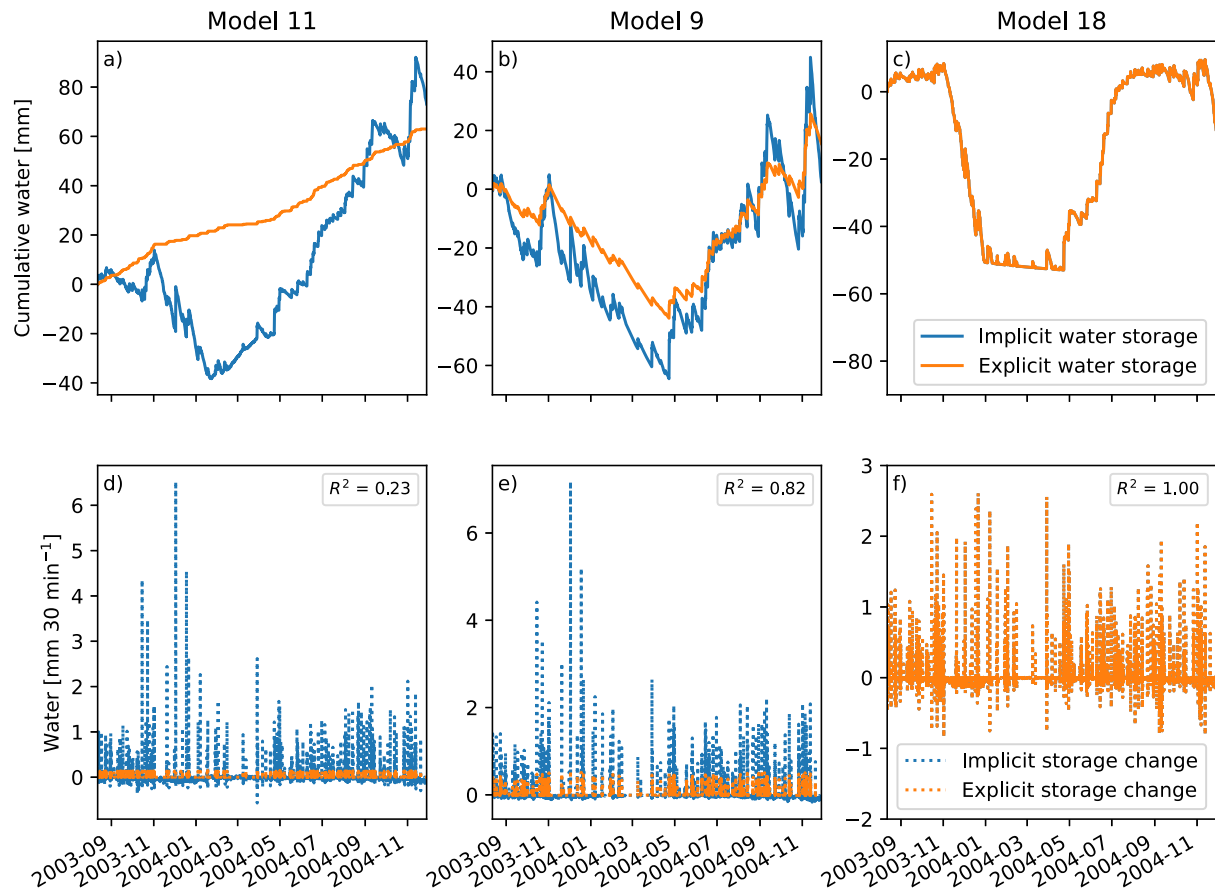


Figure 9. Illustration of (a–c) the hourly explicit (Equation 2) and implicit water storage (Equation 1) and (d–f) their derivatives both for 475 days at AU-Preston for three models with increasing coefficient of determination (R^2) of the explicit and implicit water storage change determined at (half-)hourly resolution.

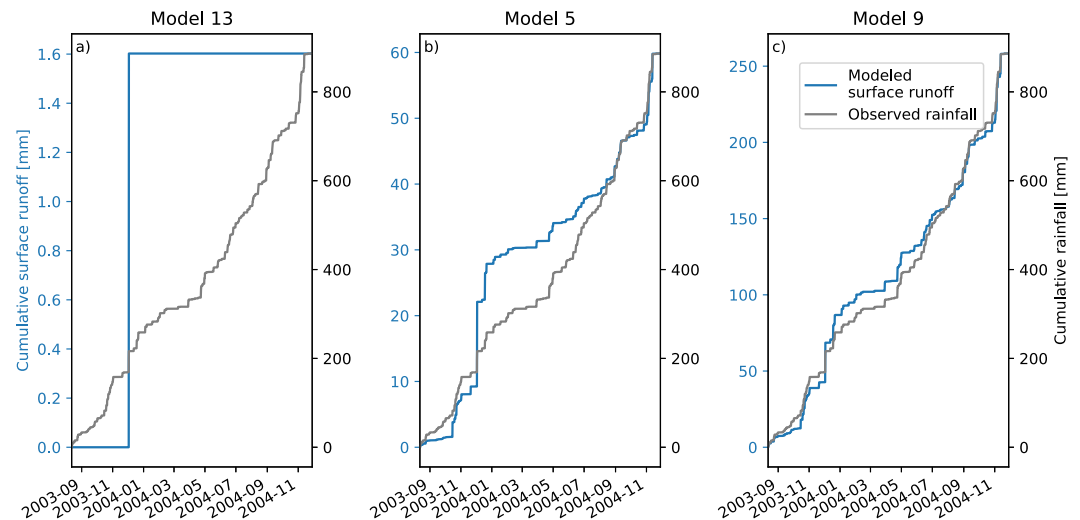


Figure 10. Illustration of surface runoff triggered for different AU-Preston precipitation events by three models (a) 13, (b) 5, and (c) 9. Note, the left-hand Y axis (surface runoff) increases (a→c), whereas the right-hand side Y axis (precipitation) is the same for all.

in more accurate latent heat fluxes. However, it is clear that the urban water balance is imperfectly incorporated into ULSMs and more proper physically based representations are required.

Five models close the water balance at all sites (Models 6, 13, 15, 18, and 19), while three never reach closure (Models 1, 3, and 5). The other models close the water balance at some sites. For several non-closing models, we identify the causes. One model implicitly assumes an infinite source or sink of soil moisture by adapting the modeled soil moisture when it exceeds hard-coded limits adding or removing water to remain within these limits (Model 11). Two other models do not fully couple all processes, such as runoff and evaporation calculations occurring without water availability feedback between processes (Models 1 and 5). Such uncoupled processes may also explain inconsistent water storage dynamics. Three models pass their internal water balance closure check but do not provide the modeled groundwater flux in the model output (Models 8, 16, and 17). We call on the modeling community to include all fluxes required to diagnose water balance closure in the model output. Three models without a snow module disregarded all snowfall creating a mismatch between real and modeled input (Model 2, 7, and 12). For one model, we suspect a very shallow soil layer causes large numerical errors resulting in an unclosed water balance (Model 4). Fortunately, model improvements should be able to eliminate these issues for most models.

Evidence is found that the models would benefit from reevaluating their runoff parameterizations. The runoff volumes are poorly captured, resulting in $I_{R,m}$ having the poorest overall pass rate (Figure 4). Runoff has not been evaluated in previous ULSM comparisons and suffers here from a lack of direct observations and small areas being modeled ($<1 \text{ km}^2$). The lack of correlation between modeled cumulative R_s and the impervious fraction is worrying given the well-documented relation (Jacobson, 2011; Shuster et al., 2005). However, many models use relatively simple approaches, such as a constant fraction of rainfall that runs off independent of site characteristics, rainfall intensity, or soil moisture state. Others use poorly constrained parameters, such as how much water is routed between sub-grid tiles. Future work could help to constrain such parameters, while the simple approaches could be improved relatively straightforwardly.

Despite the lack of evidence showing a link between the UWBR score and Q_E performance, the incomplete representation of the water balance may contribute to the poor latent heat flux performance of the ULSMs. The design of the UWBR score may not be successful in revealing an existing link between the UWBR score and Q_E performance, as the UWBR score indicators assess the water balance based on physical realism and expectations derived from the literature. While a higher UWBR score indicates a more physically consistent water balance, it may still be an incorrect simulation. The opposite is also true, as, without physical constraints, machine learning approaches show good results for Q_E (Vulova et al., 2021). Apart from that, a potential link between the water balance representation and the Q_E performance may be hidden by other elements affecting Q_E performance. These elements could be other components of the model (e.g., the energy balance representation) or human errors (e.g., erroneous parameters, assuming northern-hemisphere vegetation, and results reported in wrong units). Yet, we do find a poor performance for Q_E consistent with the literature showing Q_E is among the most challenging fluxes to model (Grimmond et al., 2011; Lipson et al., 2024). As the energy and water balance are directly connected, we hypothesize potential errors in the water balance are causing, and not being caused by, the poor performance of Q_E , as the short runoff timescales in urban areas on a neighborhood scale dictate the water availability for Q_E and not the other way around. Hence, good model performance for the latent and sensible heat flux cannot be achieved without properly representing both balances. Thus, we believe an improved representation of the water balance will assist in latent heat flux simulation and other energy fluxes.

This first systematic analysis of urban water balance modeling is an opportunistic study taking advantage of model outputs, model characterizations, and observations gathered for the Urban PLUMBER project (Lipson et al., 2022b, 2024). The Urban-PLUMBER setup affects this study via (a) the diversity of model outputs linked to their range of modeling approaches, and (b) a lack of observations for all the water balance terms. Intentionally, a wide range of modeling approaches are analyzed with both default parameters and provided parameters implemented by modelers (Lipson et al., 2024), impacting the model results and performance. For example, numerical discretization of soil layers can cause a flawed, reduced moisture drydown linked to irregular soil layer depths that enhance evaporation (MacKay et al., 2022). Ongoing land surface model developments to capture and link more processes increase both their scope and complexity, but the number of differing aspects complicates a systematic analysis aiming to attribute performance to certain aspects (Blyth et al., 2021; Fisher & Koven, 2020). To minimize human error, Urban-PLUMBER allowed resubmission of model outputs after web-based and manual

checks. As these checks did not address the water balance, we provided an additional basic analysis of the water balance results to catch other human errors with encouragement to resubmit updated outputs. Unfortunately, resubmission reduces but does not eliminate human errors. All differences other than the water balance representation hinder the attribution of the model performance to the water balance concept as they explain the large variety in model performance amongst models that capture the water balance equally accurately. Ideally, these differences would be eliminated by developing a multi-model framework in the future (Sadeq et al., 2019) and characterizing model types based on water balance approaches. Such a characterization could allow for teasing out more detailed strengths and weaknesses of water balance representations.

Lack of observations (e.g., runoff, soil moisture) prevents direct assessment for many water balance terms. These observations are challenging as both energy and water balance closure need to be considered, so observations need to cover a relatively large uniform area that also constrains the natural and anthropogenic water flows (Grimmond & Oke, 1986, 1991). A large uniform area is needed as eddy-covariance footprints vary continuously (Feigenwinter et al., 2012; Grimmond & Oke, 1991), while catchment boundaries are static. Hence, we develop a new alternative using quantitative indicators. Each indicator addresses a water balance process and checks whether it complies with physical limits, the model itself, or previous research. We refrain from weighting the indicators to minimize the score subjectivity and prevent one indicator from controlling the outcome. The systematic removal of one of the seven indicators allows us to confirm the UWBR score is not driven by one indicator.

Here, we show ULSMs produce a wide range of water balance results but often do not realistically represent important hydrological processes. Output reporting errors may cause part of the low performance. Although our results are for offline ULSMs, we expect the identified issues will persist in a coupled setting on any scale (e.g., with mesoscale and global models). ULSMs could be improved by ensuring they close the water balance and updating runoff parameterizations. Ideally, future energy-water-carbon studies will try to gather both a wider range of observations but also modeled processes. This will aid improvement of model processes and their feedbacks. However, the complexity of the urban landscape (e.g., different definitions between eddy covariance footprints, and runoff catchments) will require nested model runs and observations to ensure consistency of all. We recommend routine assessment of water balance closure in ULSM development phase applying the indicators of the UWBR score. In a broader context, both model evaluations and comparisons should extend beyond the target variables of the model to all processes that directly influence these variables. This will benefit the broader delivery of integrated urban services (WMO, 2019) and facilitate urban resilience across time scales.

Data Availability Statement

All observation data from this study are openly available at Zenodo via <https://doi.org/10.5281/zenodo.6590886> (Lipson et al., 2022a). Model results and benchmarks (Lipson & Best, 2022) for AU-Preston are archived at Zenodo. Model results for the other sites are visualized at <https://urban-plumber.github.io/sites> and will be published together with Urban-PLUMBER Phase 2.

References

- Berne, A., Delrieu, G., Creutin, J.-D., & Obled, C. (2004). Temporal and spatial resolution of rainfall measurements required for urban hydrology. *Journal of Hydrology*, 299(3–4), 166–179. [https://doi.org/10.1016/S0022-1694\(04\)00363-4](https://doi.org/10.1016/S0022-1694(04)00363-4)
- Berthier, E., Andrieu, H., & Rodríguez, F. (1999). The Rezé urban catchments database. *Water Resources Research*, 35(6), 1915–1919. <https://doi.org/10.1029/1999wr900053>
- Best, M. J., Abramowitz, G., Johnson, H., Pitman, A., Balsamo, G., Boone, A., et al. (2015). The plumbing of land surface models: Benchmarking model performance. *Journal of Hydrometeorology*, 16(3), 1425–1442. <https://doi.org/10.1175/jhm-d-14-0158.1>
- Best, M. J., & Grimmond, C. S. B. (2016). Modeling the partitioning of turbulent fluxes at urban sites with varying vegetation cover. *Journal of Hydrometeorology*, 17(10), 2537–2553. <https://doi.org/10.1175/jhm-d-15-0126.1>
- Best, M. J., Pryor, M., Clark, D., Rooney, G., Essery, R., Ménard, C., et al. (2011). The Joint UK Land Environment Simulator (JULES), model description—Part 1: Energy and water fluxes. *Geoscientific Model Development*, 4(3), 677–699. <https://doi.org/10.5194/gmd-4-677-2011>
- Björkegren, A., Grimmond, C. S. B., Kotthaus, S., & Malamud, B. (2015). CO₂ emission estimation in the urban environment: Measurement of the CO₂ storage term. *Atmospheric Environment*, 122, 775–790. <https://doi.org/10.1016/j.atmosenv.2015.10.012>
- Blyth, E. M., Arora, V. K., Clark, D. B., Dadson, S. J., De Kauwe, M. G., Lawrence, D. M., et al. (2021). Advances in land surface modelling. *Current Climate Change Reports*, 7(2), 45–71. <https://doi.org/10.1007/s40641-021-00171-5>
- Botzen, W., Martinijs, M., Bröde, P., Folkerts, M., Ignjacevic, P., Estrada, F., et al. (2020). Economic valuation of climate change-induced mortality: Age-dependent cold and heat mortality in The Netherlands. *Climatic Change*, 162(2), 545–562. <https://doi.org/10.1007/s10584-020-02797-0>
- Boussetta, S., Balsamo, G., Arduini, G., Dutra, E., McNorton, J., Choulga, M., et al. (2021). ECLand: The ECMWF land surface modelling system. *Atmosphere*, 12(6), 723. <https://doi.org/10.3390/atmos12060723>

Acknowledgments

We acknowledge the Urban-PLUMBER project team and all observation and modeling participants providing the data set for this research. We would like to thank Judith Boeke, Andrew Frost, and Valentina Marchionni for the fruitful discussions. We want to express our appreciation to the three anonymous reviewers who took the time and effort to review and help improve the manuscript. Harro Jongen acknowledges this research was supported by the WIMEK PhD Grant 2020. Mathew Lipson acknowledges support from the Australian Research Council (ARC) Centre of Excellence for Climate System Science (Grant CE110001028), National Computational Infrastructure (NCI) Australia and the Bureau of Meteorology, Australia. Gert-Jan Steeneveld acknowledges support from the Amsterdam Institute for Advanced Metropolitan Solutions (AMS Institute, project VIR16002) and the Netherlands Organization for Scientific Research (NWO, Project 864.14.007). Sue Grimmond acknowledges support from ERC Urbisphere (Grant 855055). Matthias Demuzere was supported by the ENLIGHT project, funded by the German Research Foundation (DFG) under Grant number 437467569. Ting Sun is supported by UKRI NERC Independent Research Fellowship (NE/P018637/2). Ruidong Li is supported by CSC scholarship. Keith Oleson's contribution is based upon work supported by the NSF National Center for Atmospheric Research, which is a major facility sponsored by the U.S. National Science Foundation under Cooperative Agreement No. 1852977. Chenghao Wang acknowledges support from the National Science Foundation (NSF) under Grants numbers OIA-2327435 and CNS-2301858 and the National Oceanic and Atmospheric Administration (NOAA) under Grant number NA21OAR4590361.

- Brakensiek, D. L., & Rawls, W. J. (1983). Green-Ampt infiltration model parameters for hydrologic classification of soils. In *Advances in irrigation and drainage: Surviving external pressures* (pp. 226–233).
- Bringfelt, B. (1986). Test of a forest evapotranspiration model. *SMHI*.
- Carlyle-Moses, D. E., Livesley, S., Baptista, M. D., Thom, J., & Szota, C. (2020). *Urban trees as green infrastructure for stormwater mitigation and use* (pp. 397–432). Forest-Water Interactions.
- Chakravarti, I., Laha, R., & Roy, J. (1967). *Handbook of methods of applied statistics. Volume I: Techniques of computation, descriptive methods, and statistical inference*. John Wiley and Sons, Incorporated.
- Châtelet, E. (1740). *Institutions de physique*. Paris.
- Chow, W. T. (2017). *Eddy covariance data measured at the CAP LTER flux tower located in the west Phoenix, AZ neighbourhood of Maryvale from 2011-12-16 through 2012-12-31*. Environmental Data Initiative.
- Chow, W. T., Volo, T. J., Vivoni, E. R., Jenerette, G. D., & Ruddell, B. L. (2014). Seasonal dynamics of a suburban energy balance in Phoenix, Arizona. *International Journal of Climatology*, 34(15), 3863–3880. <https://doi.org/10.1002/joc.3947>
- Christen, A., Coops, N., Crawford, B., Kellett, R., Liss, K., Olchovski, I., et al. (2011). Validation of modelled carbon-dioxide emissions from an urban neighbourhood with direct eddy-covariance measurements. *Atmospheric Environment*, 45(33), 6057–6069. <https://doi.org/10.1016/j.atmosenv.2011.07.040>
- Coutts, A. M., Beringer, J., & Tapper, N. J. (2007a). Characteristics influencing the variability of urban CO₂ fluxes in Melbourne, Australia. *Atmospheric Environment*, 41(1), 51–62. <https://doi.org/10.1016/j.atmosenv.2006.08.030>
- Coutts, A. M., Beringer, J., & Tapper, N. J. (2007b). Impact of increasing urban density on local climate: Spatial and temporal variations in the surface energy balance in Melbourne, Australia. *Journal of Applied Meteorology and Climatology*, 46(4), 477–493. <https://doi.org/10.1175/jam2462.1>
- Crawford, B., & Christen, A. (2015). Spatial source attribution of measured urban eddy covariance CO₂ fluxes. *Theoretical and Applied Climatology*, 119(3–4), 733–755. <https://doi.org/10.1007/s00704-014-1124-0>
- Crawford, B., Grimmond, C., & Christen, A. (2011). Five years of carbon dioxide fluxes measurements in a highly vegetated suburban area. *Atmospheric Environment*, 45(4), 896–905. <https://doi.org/10.1016/j.atmosenv.2010.11.017>
- Cronshey, R., Roberts, R., & Miller, N. (1985). Urban hydrology for small watersheds (TR-55 rev.). In *Hydraulics and hydrology in the small computer age* (pp. 1268–1273).
- Demuzere, M., Kittner, J., Martilli, A., Mills, G., Moede, C., Stewart, I. D., et al. (2022). A global map of local climate zones to support earth system modelling and urban scale environmental science. *Earth System Science Data Discussions*, 2022(8), 1–57. <https://doi.org/10.5194/essd-14-3835-2022>
- Ek, M., Mitchell, K., Lin, Y., Rogers, E., Grunmann, P., Koren, V., et al. (2003). Implementation of Noah land surface model advances in the National Centre's for Environmental Prediction operational mesoscale Eta model. *Journal of Geophysical Research*, 108(D22). <https://doi.org/10.1029/2002jd003296>
- Feigenwinter, C., Vogt, R., & Christen, A. (2012). Eddy covariance measurements over urban areas. Eddy covariance: A practical guide to measurement and data analysis (pp. 377–397).
- Fisher, R. A., & Koven, C. D. (2020). Perspectives on the future of land surface models and the challenges of representing complex terrestrial systems. *Journal of Advances in Modeling Earth Systems*, 12(4), e2018MS001453. <https://doi.org/10.1029/2018ms001453>
- Fletcher, T. D., Andrieu, H., & Hamel, P. (2013). Understanding, management and modelling of urban hydrology and its consequences for receiving waters: A state of the art. *Advances in Water Resources*, 51, 261–279. <https://doi.org/10.1016/j.advwatres.2012.09.001>
- Foken, T., Leuning, R., Oncley, S. R., Mauder, M., & Aubinet, M. (2012). Corrections and data quality control. In *Eddy covariance* (pp. 85–131). Springer.
- Fortuniak, K. (2003). A slab surface energy balance model (SUEB) and its application to the study on the role of roughness length in forming an urban heat island. *Acta Universitatis Wratislaviensis*, 2542, 368–377.
- Fortuniak, K., Klysik, K., & Siedlecki, M. (2006). New measurements of the energy balance components in Łódź. In *Preprints, sixth international conference on urban climate* (pp. 12–16). Göteborg.
- Fortuniak, K., Pawlak, W., & Siedlecki, M. (2013). Integral turbulence statistics over a central European city centre. *Boundary-Layer Meteorology*, 146(2), 257–276. <https://doi.org/10.1007/s10546-012-9762-1>
- Franssen, H. H., Stöckli, R., Lehner, I., Rotenberg, E., & Seneviratne, S. I. (2010). Energy balance closure of eddy-covariance data: A multisite analysis for European fluxnet stations. *Agricultural and Forest Meteorology*, 150(12), 1553–1567. <https://doi.org/10.1016/j.agrformet.2010.08.005>
- Gasparri, A., Guo, Y., Sera, F., Vicedo-Cabrera, A. M., Huber, V., Tong, S., et al. (2017). Projections of temperature-related excess mortality under climate change scenarios. *The Lancet Planetary Health*, 1(9), e360–e367. [https://doi.org/10.1016/s2542-5196\(17\)30156-0](https://doi.org/10.1016/s2542-5196(17)30156-0)
- Goret, M., Masson, V., Schoetter, R., & Moine, M.-P. (2019). Inclusion of CO₂ flux modelling in an urban canopy layer model and an evaluation over an old European city centre. *Atmospheric Environment X*, 3, 100042. <https://doi.org/10.1016/j.aeaoa.2019.100042>
- Grimmond, C. S. B. (2006). Progress in measuring and observing the urban atmosphere. *Theoretical and Applied Climatology*, 84(1), 3–22. <https://doi.org/10.1007/s00704-005-0140-5>
- Grimmond, C. S. B., Best, M. J., Barlow, J., Arnfield, A., Baik, J.-J., Baklanov, A., et al. (2009). Urban surface energy balance models: Model characteristics and methodology for a comparison study. In *Meteorological and air quality models for urban areas* (pp. 97–123). Springer.
- Grimmond, C. S. B., Blackett, M., Best, M. J., Baik, J.-J., Belcher, S., Beringer, J., et al. (2011). Initial results from phase 2 of the international urban energy balance model comparison. *International Journal of Climatology*, 31(2), 244–272. <https://doi.org/10.1002/joc.2227>
- Grimmond, C. S. B., Blackett, M., Best, M. J., Barlow, J., Baik, J., Belcher, S., et al. (2010). The international urban energy balance models comparison project: First results from phase 1. *Journal of Applied Meteorology and Climatology*, 49(6), 1268–1292. <https://doi.org/10.1175/2010jame2354.1>
- Grimmond, C. S. B., & Oke, T. R. (1986). Urban water balance: 2. Results from a suburb of Vancouver, British Columbia. *Water Resources Research*, 22(10), 1404–1412. <https://doi.org/10.1029/wr022i010p01404>
- Grimmond, C. S. B., & Oke, T. R. (1991). An evapotranspiration-interception model for urban areas. *Water Resources Research*, 27(7), 1739–1755. <https://doi.org/10.1029/91wr00557>
- Grimmond, C. S. B., & Oke, T. R. (2002). Turbulent heat fluxes in urban areas: Observations and a local-scale urban meteorological parameterization scheme (LUMPS). *Journal of Applied Meteorology and Climatology*, 41(7), 792–810. [https://doi.org/10.1175/1520-0450\(2002\)041<0792:thfua>2.0.co;2](https://doi.org/10.1175/1520-0450(2002)041<0792:thfua>2.0.co;2)
- Hamdi, R., Kusaka, H., Doan, Q.-V., Cai, P., He, H., Luo, G., et al. (2020). The state-of-the-art of urban climate change modeling and observations. *Earth Systems and Environment*, 4(4), 1–16. <https://doi.org/10.1007/s41748-020-00193-3>

- Hamdi, R., & Schayes, G. (2007). Validation of Martilli's urban boundary layer scheme with measurements from two mid-latitude European cities. *Atmospheric Chemistry and Physics*, 7(17), 4513–4526. <https://doi.org/10.5194/acp-7-4513-2007>
- Heaviside, C., Vardoulakis, S., & Cai, X.-M. (2016). Attribution of mortality to the urban heat island during heatwaves in the West Midlands, UK. *Environmental Health*, 15(1), 49–59. <https://doi.org/10.1186/s12940-016-0100-9>
- Hellsten, A., Luukkonen, S.-M., Steinfeld, G., Kanani-Sühring, F., Markkanen, T., Järvi, L., et al. (2015). Footprint evaluation for flux and concentration measurements for an urban-like canopy with coupled Lagrangian stochastic and large-eddy simulation models. *Boundary-Layer Meteorology*, 157(2), 191–217. <https://doi.org/10.1007/s10564-015-0062-4>
- Henderson-Sellers, A., McGuffie, K., & Pitman, A. (1996). The project for intercomparison of land-surface parametrization schemes (PILPS): 1992 to 1995. *Climate Dynamics*, 12(12), 849–859. <https://doi.org/10.1007/s003820050147>
- Hengl, T. (2018). Sand content in % (kg/kg) at 6 standard depths (0, 10, 30, 60, 100 and 200 cm) at 250 m resolution (version v0. 2) [Dataset]. *Zenodo*. <https://doi.org/10.5281/zenodo.2525662>
- Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., et al. (2020). The ERA5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*, 146(730), 1999–2049. <https://doi.org/10.1002/qj.3803>
- Hertwig, D., Grimmond, C. S. B., Hendry, M. A., Saunders, B., Wang, Z., Jeoffrion, M., et al. (2020). Urban signals in high-resolution weather and climate simulations: Role of urban land-surface characterisation. *Theoretical and Applied Climatology*, 142(1), 701–728. <https://doi.org/10.1007/s00704-020-03294-1>
- Hirano, T., Sugawara, H., Murayama, S., & Kondo, H. (2015). Diurnal variation of CO₂ flux in an urban area of Tokyo. *Sola*, 11(0), 100–103. <https://doi.org/10.2151/sola.2015-024>
- Hirschi, M., Michel, D., Lehner, I., & Seneviratne, S. I. (2017). A site-level comparison of lysimeter and eddy covariance flux measurements of evapotranspiration. *Hydrology and Earth System Sciences*, 21(3), 1809–1825. <https://doi.org/10.5194/hess-21-1809-2017>
- Hong, J.-W., Hong, J., Chun, J., Lee, Y. H., Chang, L.-S., Lee, J.-B., et al. (2019). Comparative assessment of net CO₂ exchange across an urbanization gradient in Korea based on eddy covariance measurements. *Carbon Balance and Management*, 14(1), 1–18. <https://doi.org/10.1186/s13021-019-0128-6>
- Hong, J.-W., Lee, K., & Hong, J. (2020). Observational data of Ochang and Jungnang in Korea [Dataset]. EAPL at Yonsei University. https://doi.org/10.22647/EAPL-OC_JN2021
- Hong, S.-O., Kim, J., Byun, Y.-H., Hong, J., Hong, J.-W., Lee, K., et al. (2023). Intra-urban variations of the CO₂ fluxes at the surface-atmosphere interface in the Seoul Metropolitan Area. *Asia-Pacific Journal of Atmospheric Sciences*, 59(4), 1–15. <https://doi.org/10.1007/s13143-023-00324-6>
- Howard, L. (1833). The climate of London deduced from meteorological observations made in the metropolis and various places around it. (Vol. 3).
- Ishidoya, S., Sugawara, H., Terao, Y., Kaneyasu, N., Aoki, N., Tsuboi, K., & Kondo, H. (2020). O₂: CO₂ exchange ratio for net turbulent flux observed in an urban area of Tokyo, Japan, and its application to an evaluation of anthropogenic CO₂ emissions. *Atmospheric Chemistry and Physics*, 20(9), 5293–5308. <https://doi.org/10.5194/acp-20-5293-2020>
- Jacobson, C. R. (2011). Identification and quantification of the hydrological impacts of imperviousness in urban catchments: A review. *Journal of Environmental Management*, 92(6), 1438–1448. <https://doi.org/10.1016/j.jenvman.2011.01.018>
- Järvi, L., Grimmond, C. S. B., & Christen, A. (2011). The Surface Urban Energy and Water balance Scheme (SUEWS): Evaluation in Los Angeles and Vancouver. *Journal of Hydrology*, 411(3–4), 219–237. <https://doi.org/10.1016/j.jhydrol.2011.10.001>
- Järvi, L., Rannik, Ü., Kokkonen, T. V., Kurppa, M., Karppinen, A., Kouznetsov, R. D., et al. (2018). Uncertainty of eddy covariance flux measurements over an urban area based on two towers. *Atmospheric Measurement Techniques*, 11(10), 5421–5438. <https://doi.org/10.5194/amt-11-5421-2018>
- Jongen, H. J., Steeneveld, G.-J., Beringer, J., Christen, A., Chrysoulakis, N., Fortuniak, K., et al. (2022). Urban water storage capacity inferred from observed evapotranspiration recession. *Geophysical Research Letters*, 49(3), e2021GL096069. <https://doi.org/10.1029/2021gl096069>
- Karsisto, P., Fortelius, C., Demuzere, M., Grimmond, C. S. B., Oleson, K., Kouznetsov, R., et al. (2016). Seasonal surface urban energy balance and wintertime stability simulated using three land-surface models in the high-latitude city Helsinki. *Quarterly Journal of the Royal Meteorological Society*, 142(694), 401–417. <https://doi.org/10.1002/qj.2659>
- Klaassen, W., Bosveld, F., & De Water, E. (1998). Water storage and evaporation as constituents of rainfall interception. *Journal of Hydrology*, 212, 36–50. [https://doi.org/10.1016/s0022-1694\(98\)00200-5](https://doi.org/10.1016/s0022-1694(98)00200-5)
- Kokkonen, T., Grimmond, C. S. B., Christen, A., Oke, T., & Järvi, L. (2018). Changes to the water balance over a century of urban development in two neighborhoods: Vancouver, Canada. *Water Resources Research*, 54(9), 6625–6642. <https://doi.org/10.1029/2017wr022445>
- Koopmans, S., Heusinkveld, B., & Steeneveld, G. (2020). A standardized Physical Equivalent Temperature urban heat map at 1-m spatial resolution to facilitate climate stress tests in The Netherlands. *Building and Environment*, 181, 106984. <https://doi.org/10.1016/j.buildenv.2020.106984>
- Kotthaus, S., & Grimmond, C. S. B. (2014a). Energy exchange in a dense urban environment—Part II: Impact of spatial heterogeneity of the surface. *Urban Climate*, 10, 281–307. <https://doi.org/10.1016/j.uclim.2013.10.001>
- Kotthaus, S., & Grimmond, C. S. B. (2014b). Energy exchange in a dense urban environment—Part I: Temporal variability of long-term observations in central London. *Urban Climate*, 10, 261–280. <https://doi.org/10.1016/j.uclim.2013.10.002>
- Kowalczyk, E., Wang, Y., Law, R., Davies, H., McGregor, J., & Abramowitz, G. (2006). The CSIRO Atmosphere Biosphere Land Exchange (CABLE) model for use in climate models and as an offline model. *CSIRO Marine and Atmospheric Research Paper*, 13, 42.
- Krayenhoff, E. S., & Voogt, J. A. (2007). A microscale three-dimensional urban energy balance model for studying surface temperatures. *Boundary-Layer Meteorology*, 123(3), 433–461. <https://doi.org/10.1007/s10564-006-9153-6>
- Kusaka, H., Kondo, H., Kikegawa, Y., & Kimura, F. (2001). A simple single-layer urban canopy model for atmospheric models: Comparison with multi-layer and slab models. *Boundary-Layer Meteorology*, 101(3), 329–358. <https://doi.org/10.1023/a:1019207923078>
- Lavoisier, A. L. (1789). *Traite elementaire de chimie*.
- Lemonsu, A., Viguie, V., Daniel, M., & Masson, V. (2015). Vulnerability to heat waves: Impact of urban expansion scenarios on urban heat island and heat stress in Paris (France). *Urban Climate*, 14, 586–605. <https://doi.org/10.1016/j.uclim.2015.10.007>
- Leopold, L. B. (1968). *Hydrology for urban land planning: A guidebook on the hydrologic effects of urban land use* (Vol. 554). US Geological Survey. <https://doi.org/10.3133/cir554>
- Li, D., Liao, W., Rigden, A. J., Liu, X., Wang, D., Malyshev, S., & Shevliakova, E. (2019). Urban heat island: Aerodynamics or imperviousness? *Science Advances*, 5(4), eaau4299. <https://doi.org/10.1126/sciadv.aau4299>
- Lipson, M. J., & Best, M. (2022). Benchmarks for the Urban-PLUMBER model evaluation project Phase 1 (AU-Preston) (Version v1) [Dataset]. *Zenodo*. <https://doi.org/10.5281/zenodo.7330052>

- Lipson, M. J., Grimmond, C. S. B., Best, M., Abramowitz, G., Coutts, A., Tapper, N., et al. (2024). Evaluation of 30 urban land surface models in the Urban-PLUMBER project: Phase 1 results. *Quarterly Journal of the Royal Meteorological Society*, 150(758), 126–169. <https://doi.org/10.1002/qj.4589>
- Lipson, M. J., Grimmond, C. S. B., Best, M., Chow, W. T., Christen, A., Chrysoulakis, N., et al. (2022a). Site data archive for “Harmonized gap-filled dataset from 20 urban flux tower sites” for the Urban-PLUMBER project (Version v1) [Dataset]. *Zenodo*. <https://doi.org/10.5281/zenodo.7104984>
- Lipson, M. J., Grimmond, C. S. B., Best, M. J., Chow, W. T., Christen, A., Chrysoulakis, N., et al. (2022b). *Harmonized gap-filled datasets from 20 urban flux tower sites* (pp. 1–29). Earth System Science Data Discussions.
- Lipson, M. J., Grimmond, S., Best, M., Abramowitz, G., Coutts, A., Tapper, N., et al. (2023). *The Urban-PLUMBER model evaluation project: Phase 1 results*. 11th International Conference on Urban Climate.
- Lipson, M. J., Thatcher, M., Hart, M. A., & Pitman, A. (2018). A building energy demand and urban land surface model. *Quarterly Journal of the Royal Meteorological Society*, 144(714), 1572–1590. <https://doi.org/10.1002/qj.3317>
- MacKay, M. D., Meyer, G., & Melton, J. R. (2022). On the discretization of Richards equation in Canadian land surface models. *Atmosphere-Ocean*, 61, 1–11. <https://doi.org/10.1080/07055900.2022.2096558>
- Manabe, S. (1969). Climate and the ocean circulation: I. The atmospheric circulation and the hydrology of the earth's surface. *Monthly Weather Review*, 97(11), 739–774. [https://doi.org/10.1175/1520-0493\(1969\)097<0739:catoc>2.3.co;2](https://doi.org/10.1175/1520-0493(1969)097<0739:catoc>2.3.co;2)
- Masson, V., Gomes, L., Pigeon, G., Lioussse, C., Pont, V., Lagouarde, J.-P., et al. (2008). The Canopy and Aerosol Particles Interactions in Toulouse Urban Layer (CAPITOL) experiment. *Meteorology and Atmospheric Physics*, 102(3–4), 135–157. <https://doi.org/10.1007/s00703-008-0289-4>
- Mauder, M., Foken, T., & Cuxart, J. (2020). Surface-energy-balance closure over land: A review. *Boundary-Layer Meteorology*, 177(2), 395–426. <https://doi.org/10.1007/s10546-020-00529-6>
- McNorton, J., Arduini, G., Bousset, N., Agustí-Panareda, A., Balsamo, G., Boussetta, S., et al. (2021). An urban scheme for the ECMWF integrated forecasting system: Single-column and global offline application. *Journal of Advances in Modeling Earth Systems*, 13(6), e2020MS002375. <https://doi.org/10.1029/2020ms002375>
- Meili, N., Manoli, G., Burlando, P., Bou-Zeid, E., Chow, W. T., Coutts, A. M., et al. (2020). An urban ecohydrological model to quantify the effect of vegetation on urban climate and hydrology (UT&C v1. 0). *Geoscientific Model Development*, 13(1), 335–362. <https://doi.org/10.5194/gmd-13-335-2020>
- Menzer, O., & McFadden, J. P. (2017). Statistical partitioning of a three-year time series of direct urban net CO₂ flux measurements into biogenic and anthropogenic components. *Atmospheric Environment*, 170, 319–333. <https://doi.org/10.1016/j.atmosenv.2017.09.049>
- Mitchell, V. G., Mein, R. G., & McMahon, T. A. (2001). Modelling the urban water cycle. *Environmental Modelling and Software*, 16(7), 615–629. [https://doi.org/10.1016/s1364-8152\(01\)00029-9](https://doi.org/10.1016/s1364-8152(01)00029-9)
- Morin, E., Enzel, Y., Shamir, U., & Garti, R. (2001). The characteristic time scale for basin hydrological response using radar data. *Journal of Hydrology*, 252(1–4), 85–99. [https://doi.org/10.1016/s0022-1694\(01\)00451-6](https://doi.org/10.1016/s0022-1694(01)00451-6)
- Nachtergaele, F. (2001). Soil taxonomy—A basic system of soil classification for making and interpreting soil surveys. *Geoderma*, 99(3–4), 336–337. [https://doi.org/10.1016/s0016-7061\(00\)00097-5](https://doi.org/10.1016/s0016-7061(00)00097-5)
- Nordbo, A., Järvi, L., Haapanala, S., Moilanen, J., & Vesala, T. (2013). Intra-city variation in urban morphology and turbulence structure in Helsinki, Finland. *Boundary-Layer Meteorology*, 146(3), 469–496. <https://doi.org/10.1007/s10546-012-9773-y>
- Oke, T. R. (1982). The energetic basis of the urban heat island. *Quarterly Journal of the Royal Meteorological Society*, 108(455), 1–24. <https://doi.org/10.1002/qj.49710845502>
- Oleson, K. W., Bonan, G. B., Feddema, J., Vertenstein, M., & Grimmond, C. S. B. (2008). An urban parameterization for a global climate model. Part I: Formulation and evaluation for two cities. *Journal of Applied Meteorology and Climatology*, 47(4), 1038–1060. <https://doi.org/10.1175/2007jamc1597.1>
- Oleson, K. W., & Feddema, J. (2020). Parameterization and surface data improvements and new capabilities for the Community Land Model Urban (CLMU). *Journal of Advances in Modeling Earth Systems*, 12(2), e2018MS001586. <https://doi.org/10.1029/2018ms001586>
- Pawlak, W., Fortuniak, K., & Siedlecki, M. (2011). Carbon dioxide flux in the centre of Łódź, Poland — Analysis of a 2-year eddy covariance measurement data set. *International Journal of Climatology*, 31(2), 232–243. <https://doi.org/10.1002/joc.2247>
- Peters, E. B., Hiller, R. V., & McFadden, J. P. (2011). Seasonal contributions of vegetation types to suburban evapotranspiration. *Journal of Geophysical Research*, 116(G1), G01003. <https://doi.org/10.1029/2010jg001463>
- Petrucchi, R. H., Herring, F. G., & Madura, J. D. (2010). *General chemistry: Principles and modern applications*. Pearson Prentice Hall.
- Porson, A., Clark, P. A., Harman, I., Best, M. J., & Belcher, S. (2010). Implementation of a new urban energy budget scheme in the MetUM. Part I: Description and idealized simulations. *Quarterly Journal of the Royal Meteorological Society*, 136(651), 1514–1529. <https://doi.org/10.1002/qj.668>
- Randrup, T. B., McPherson, E. G., & Costello, L. R. (2001). Tree root intrusion in sewer systems: A review of extent and costs. *Journal of Infrastructure Systems*, 7(1), 26–31. [https://doi.org/10.1061/\(asce\)1076-0342\(2001\)7:1\(26\)](https://doi.org/10.1061/(asce)1076-0342(2001)7:1(26))
- Ronda, R., Steeneveld, G., Heusinkveld, B., Attema, J., & Holtslag, A. (2017). Urban fine scale forecasting reveals weather conditions with unprecedented detail. *Bulletin of the American Meteorological Society*, 98(12), 2675–2688. <https://doi.org/10.1175/bams-d-16-0297.1>
- Ross, S. L., & Oke, T. (1988). Tests of three urban energy balance models. *Boundary-Layer Meteorology*, 44(1), 73–96. <https://doi.org/10.1007/bf00117293>
- Roth, M., Jansson, C., & Velasco, E. (2017). Multi-year energy balance and carbon dioxide fluxes over a residential neighbourhood in a tropical city. *International Journal of Climatology*, 37(5), 2679–2698. <https://doi.org/10.1002/joc.4873>
- Ryu, Y.-H., Baik, J.-J., & Lee, S.-H. (2011). A new single-layer urban canopy model for use in mesoscale atmospheric models. *Journal of Applied Meteorology and Climatology*, 50(9), 1773–1794. <https://doi.org/10.1175/2011jamc2665.1>
- Sadegh, M., AghaKouchak, A., Flores, A., Mallakpour, I., & Nikoo, M. R. (2019). A multi-model nonstationary rainfall-runoff modeling framework: Analysis and toolbox. *Water Resources Management*, 33(9), 3011–3024. <https://doi.org/10.1007/s11269-019-02283-y>
- Saxton, K., Rawls, W., Romberger, J. S., & Papendick, R. (1986). Estimating generalized soil-water characteristics from texture. *Soil Science Society of America Journal*, 50(4), 1031–1036. <https://doi.org/10.2136/sssaj1986.03615995005000040054x>
- Schulz, J.-P., & Vogel, G. (2020). Improving the processes in the land surface scheme TERRA: Bare soil evaporation and skin temperature. *Atmosphere*, 11(5), 513. <https://doi.org/10.3390/atmos11050513>
- Shuster, W. D., Bonta, J., Thurston, H., Warnemuende, E., & Smith, D. (2005). Impacts of impervious surface on watershed hydrology: A review. *Urban Water Journal*, 2(4), 263–275. <https://doi.org/10.1080/15730620500386529>

- Stagakis, S., Chrysoulakis, N., Spyridakis, N., Feigenwinter, C., & Vogt, R. (2019). Eddy covariance measurements and source partitioning of CO₂ emissions in an urban environment: Application for Heraklion, Greece. *Atmospheric Environment*, 201, 278–292. <https://doi.org/10.1016/j.atmosenv.2019.01.009>
- Steenneveld, G.-J., van der Horst, S., & Heusinkveld, B. (2020). Observing the surface radiation and energy balance, carbon dioxide and methane fluxes over the city centre of Amsterdam. In *EGU general assembly conference abstracts* (p. 1547).
- Stewart, I. D., & Oke, T. R. (2012). Local climate zones for urban temperature studies. *Bulletin of the American Meteorological Society*, 93(12), 1879–1900. <https://doi.org/10.1175/bams-d-11-00019.1>
- Sun, R., Wang, Y., & Chen, L. (2018). A distributed model for quantifying temporal-spatial patterns of anthropogenic heat based on energy consumption. *Journal of Cleaner Production*, 170, 601–609. <https://doi.org/10.1016/j.jclepro.2017.09.153>
- Templeton, N. P., Vivoni, E. R., Wang, Z.-H., & Schreiner-McGraw, A. P. (2018). Quantifying water and energy fluxes over different urban land covers in Phoenix, Arizona. *Journal of Geophysical Research: Atmospheres*, 123(4), 2111–2128. <https://doi.org/10.1002/2017jd027845>
- Tewari, M., Chen, F., Kusaka, H., & Miao, S. (2007). Coupled WRF/Unified Noah/urban-canopy modeling system. In *Ncar WRF documentation*, NCAR, Boulder (Vol. 122, pp. 1–22). Citeseer.
- Thatcher, M., & Hurley, P. (2012). Simulating Australian urban climate in a mesoscale atmospheric numerical model. *Boundary-Layer Meteorology*, 142(1), 149–175. <https://doi.org/10.1007/s10546-011-9663-8>
- Twine, T. E., Kustas, W., Norman, J., Cook, D., Houser, P., Meyers, T., et al. (2000). Correcting eddy-covariance flux underestimates over a grassland. *Agricultural and Forest Meteorology*, 103(3), 279–300. [https://doi.org/10.1016/s0168-1923\(00\)00123-4](https://doi.org/10.1016/s0168-1923(00)00123-4)
- United Nations. (2018). *World urbanization prospects, the 2018 revision*. UN Department of Economic and Social Affairs.
- Van de Vijver, E., Delbecq, N., Verdoodt, A., & Seuntjens, P. (2020). Estimating the urban soil information gap using exhaustive land cover data: The example of Flanders, Belgium. *Geoderma*, 372, 114371. <https://doi.org/10.1016/j.geoderma.2020.114371>
- Velasco, E., Perrusquia, R., Jiménez, E., Hernández, F., Camacho, P., Rodríguez, S., et al. (2014). Sources and sinks of carbon dioxide in a neighbourhood of Mexico City. *Atmospheric Environment*, 97, 226–238. <https://doi.org/10.1016/j.atmosenv.2014.08.018>
- Velasco, E., Pressley, S., Grivicke, R., Allwine, E., Molina, L. T., & Lamb, B. (2011). Energy balance in urban Mexico City: Observation and parameterization during the MILAGRO/MCMA-2006 field campaign. *Theoretical and Applied Climatology*, 103(3–4), 501–517. <https://doi.org/10.1007/s00704-010-0314-7>
- Vulova, S., Meier, F., Rocha, A. D., Quanz, J., Nouri, H., & Kleinschmit, B. (2021). Modeling urban evapotranspiration using remote sensing, flux footprints, and artificial intelligence. *Science of the Total Environment*, 786, 147293. <https://doi.org/10.1016/j.scitotenv.2021.147293>
- Wald, A. (1943). Tests of statistical hypotheses concerning several parameters when the number of observations is large. *Transactions of the American Mathematical Society*, 54(3), 426–482. <https://doi.org/10.2307/1990256>
- Walsh, C. J., Fletcher, T. D., & Ladson, A. R. (2005). Stream restoration in urban catchments through redesigning storm water systems: Looking to the catchment to save the stream. *Journal of the North American Benthological Society*, 24(3), 690–705. <https://doi.org/10.1899/04-020.1>
- Wang, C., Wang, Z.-H., & Ryu, Y.-H. (2021). A single-layer urban canopy model with transmissive radiation exchange between trees and street canyons. *Building and Environment*, 191, 107593. <https://doi.org/10.1016/j.buildenv.2021.107593>
- Wang, Y. P., Kowalczyk, E., Leuning, R., Abramowitz, G., Raupach, M. R., Pak, B., et al. (2011). Diagnosing errors in a land surface model (CABLE) in the time and frequency domains. *Journal of Geophysical Research*, 116(G1), G01034. <https://doi.org/10.1029/2010jg001385>
- Wang, Z.-H., Bou-Zeid, E., & Smith, J. A. (2013). A coupled energy transport and hydrological model for urban canopies evaluated using a wireless sensor network. *Quarterly Journal of the Royal Meteorological Society*, 139(675), 1643–1657. <https://doi.org/10.1002/qj.2032>
- Ward, H. C., Evans, J. G., & Grimmond, C. S. B. (2013). Multi-season eddy covariance observations of energy, water and carbon fluxes over a suburban area in Swindon, UK. *Atmospheric Chemistry and Physics*, 13(9), 4645–4666. <https://doi.org/10.5194/acp-13-4645-2013>
- Ward, H. C., Kotthaus, S., Järvi, L., & Grimmond, C. S. B. (2016). Surface Urban Energy and Water balance Scheme (SUEWS): Development and evaluation at two UK sites. *Urban Climate*, 18, 1–32. <https://doi.org/10.1016/j.uclim.2016.05.001>
- Wenzel, H. G., Jr., & Voorhees, M. L. (1981). *Evaluation of the urban design storm concept*. University of Illinois at Urbana-Champaign. Water Resources Center.
- Willmott, C. J. (1982). Some comments on the evaluation of model performance. *Bulletin of the American Meteorological Society*, 63(11), 1309–1313. [https://doi.org/10.1175/1520-0477\(1982\)063<1309:scoteo>2.0.co;2](https://doi.org/10.1175/1520-0477(1982)063<1309:scoteo>2.0.co;2)
- WMO. (2019). *Guidance on integrated urban hydrometeorological, climate and environmental services - volume I*. Geneva.
- Wouters, H., Demuzere, M., De Ridder, K., & van Lipzig, N. P. (2015). The impact of impervious water-storage parametrization on urban climate modelling. *Urban Climate*, 11, 24–50. <https://doi.org/10.1016/j.uclim.2014.11.005>
- Yang, Z.-L., Dickinson, R., Henderson-Sellers, A., & Pitman, A. (1995). Preliminary study of spin-up processes in land surface models with the first stage data of project for intercomparison of land surface parameterization schemes phase 1 (a). *Journal of Geophysical Research*, 100(D8), 16553–16578. <https://doi.org/10.1029/95jd01076>
- Yao, L., Wei, W., & Chen, L. (2016). How does imperviousness impact the urban rainfall-runoff process under various storm cases? *Ecological Indicators*, 60, 893–905. <https://doi.org/10.1016/j.ecolind.2015.08.041>
- Yu, M., Wu, H., Yin, J., Liang, X., & Miao, S. (2022). Improved delineation of urban hydrological processes in coupled regional climate models. *Water Resources Research*, 58(11), e2022WR032695. <https://doi.org/10.1029/2022wr032695>
- Zeisl, P., Mair, M., Kastlunger, U., Bach, P. M., Rauch, W., Sitzenfrie, R., & Kleidorfer, M. (2018). Conceptual urban water balance model for water policy testing: An approach for large scale investigation. *Sustainability*, 10(3), 716. <https://doi.org/10.3390/su10030716>
- Zhou, Q., Leng, G., Su, J., & Ren, Y. (2019). Comparison of urbanization and climate change impacts on urban flood volumes: Importance of urban planning and drainage adaptation. *Science of the Total Environment*, 658, 24–33. <https://doi.org/10.1016/j.scitotenv.2018.12.184>