

Dear Dr. Buytaert,

Thank you for offering to publish our contribution “*Technical Note: Evaluation and bias correction of an observations-based global runoff dataset using historical streamflow observations from small tropical catchments in the Philippines*” subject to revisions. We offer our point by point responses followed by a tracked changes document, as well as a clean manuscript uploaded to the system.

We viewed the reviewer comments as robust and constructive. As such, we have implemented all suggestions from the reviewers and in conjunction with you (over email) have chosen to keep the manuscript as a HESS “Technical Note”. Further, the reviewer suggestions for more robust discussion of the statistics/uncertainties, an additional figure (and new sub panels to existing figures), and qualification of the utility of GRUN and our bias correction were all implemented.

We would also like to note that in addition to the two reviews we solicited feedback from the Gionata Ghiggi the first author of the global runoff product (GRUN) used in this work. Ghiggi provided extensive comments and clarifications, all of which we have incorporated into this manuscript, we have thus acknowledged Ghiggi to recognize this contribution that has improved this work.

Correspondence associated with this manuscript should be directed to Daniel E. Ibarra and Carlos Primo C. David. **All responses are in red below.** We trust that our article is now suitable for publication in HESS. Thank you again for handling our manuscript.

Sincerely,

A handwritten signature in black ink, appearing to read "Daniel E. Ibarra". The script is fluid and cursive, with the first name being more prominent.

Daniel Ibarra

Miller Institute and UC President's Postdoctoral Fellow, Earth and Planetary Science, University of California, Berkeley

Visiting Assistant Professor of Environment and Society, Brown University

## Response to review by Jaivime Evaristo

General comments: It has been an absolute pleasure reading through this contribution. By using monthly runoff observations from 55 catchments in the Philippines with more than 10 years of data between 1946 and 2014, Ibarra et al. evaluated the possible utility (and veracity) of a recently published global runoff product GRUN\_v1 (Ghiggi et al., 2019). They showed significant albeit weak correlation between their data and GRUN model predictions, and somewhat improved model-data correspondence using volumetric efficiency (VE) and log-transformed NSE criteria. Among others, Ibarra et al. demonstrated systematic over- and underprediction of baseflow during the dry months, and underprediction of peak flow in some wet months in most catchments. To go above and beyond a simple demonstration of model-data correspondence, the authors proposed a two-step bias correction procedure that particularly addresses GRUN underpredictions during the wettest months. The authors suggested that the utility of GRUN can be extended to other ungauged tropical basins if a similar bias correction methodology is applied.

While GRUN\_v1 was trained and validated using GISM and GRDC, respectively, none of the corresponding GISM and GRDC data from the Philippines was used in GRUN\_v1. Thus, it is worth noting that this contribution by Ibarra et al. is indeed an independent test of GRUN\_V1 runoff reconstruction.

There is no doubt that the broader community will stand to benefit from Ibarra et al.'s analysis. The scientific and engineering literature on water resources continues to “suffer” from a mid- to high-latitude bias. Ibarra et al.'s work represents a substantial contribution to reducing this bias and increasing our understanding of tropical hydrology, particularly with respect to the implications of their work for the ungauged tropical basins. Moreover, I can only hope that the community will also commend Ibarra et al. for making these Philippine datasets publicly available, which may prove useful for similar and sundry purposes. These favorable comments notwithstanding, I raise some [relatively minor] points that when addressed may only serve to improve this contribution.

**We thank the reviewer for thoughtful comments and points to address. All of them will be incorporated into our revised manuscript.**

Specific comments: (1) On bias correction at the national scale: Is there any particular practical significance for the bias correction at the national level, as opposed to, say, at the basin level or per climate types? For example, the per-climate-type analysis seems to show some interesting patterns (Figs. 3&4), and so as at the basin level or catchment size (Fig. 2B). This comment of course assumes that a sufficiently wide range of flows are similarly captured at these levels of abstraction as at the national level, thereby, making log-transformations meaningful. Such seems to be the case per climate type based on the scatterplots in Figure 3. In any case, it might be useful to know why the bias correction was performed at the aggregate national level and not at [or not in addition to] sub-national levels

This was partially due to sample size constraints and because climate type does not translate into regions (see colored areas in Figure 1). Sample size is an issue particularly for Type 2 catchments, we are limited to 6 catchments, and catchment areas are not evenly distributed among the catchments shown in figure 2B (i.e., Type 1 has few small catchments and type 4 has few large catchments). Additionally, in balancing this reviewers suggestions with those of Anonymous Reviewer #2 who felt that we were overly optimistic in our statement of utility with respect to using GRUN for ungauged small tropical catchments we feel that only providing a national level bias correction is most appropriate for use at a regional/national scale.

(2) Parameter uncertainty: I would encourage the authors to also perform uncertainty estimation on their slope ( $m=0.774$ ) and intercept ( $b=0.099$ ) parameters, possibly via bootstrapping. This would make their proposed bias correction method more robust and bounded. Suggested references follow:

Efron, B. (1981), “Nonparametric Standard Errors and Confidence Intervals,” *The Canadian Journal of Statistics* 9:2, 139–158.

Rubin, D. (1981), “The Bayesian Bootstrap”, *The Annals of Statistics*, 9:1, 130–134

This is an excellent suggestion, we have performed a bootstrap analysis on our regression and report the 68% (1 sigma) confidence intervals for bias correction, the text will now read: “Thus, we first add the mean  $\log_{10}(\text{runoff})$  difference between the observations and the predicted values ( $0.117 \pm 0.022$ ). Following this, using the `lm` function in R, we fit a linear regression between the observations and the GRUN predicted values ( $\log_{10}(\text{runoff, observed}) = m \times 270 \log_{10}(\text{runoff, predicted}) + b$ ) and correct the predicted values using the slope ( $m=0.774 \pm 0.025$ ) and intercept ( $b=0.099 \pm 0.006$ ) derived from this regression. Uncertainties reported here are 68% confidence intervals and were assessed by bootstrap resampling 1,000 observation and prediction pairs without replacement 10,000 times.”

(3) On transformation bias in curve fitting: One utility of Ibarra et al.’s work is on possibly using GRUN for other ungauged basins in the tropics and applying a similar bias correction as proposed (L301). Because these corrections are in log-log space, the user may then need to back-transform (or antilog) to obtain the “corrected” runoff values. Such toggling has long been shown to carry some inherent statistical bias that [also] needs to be corrected, as succinctly discussed by Ferguson (1986). This bias can be non-trivial and results from the use of least squares regression in estimating the logarithms of, say, runoff in ungauged basins. Without repeating here the arguments that Ferguson most effectively articulated in 1986 (and Miller 1984), the authors may find it worthwhile to reflect on the implications of this [possible] statistical bias on their proposed method for bias correction.

Ferguson, R. I. River Loads Underestimated by Rating Curves. *Water Resour. Res.* 22, 74–76 (1986).

Miller, D. M. Reducing Transformation Bias in Curve Fitting. *Am. Stat.* 38, 124 (1984).

We agree with this comment entirely and thank the reviewer for pointing out these references. This is an important point, following the calculation of the bias correction values in our discussion we will add several sentences to discuss this potential bias and note that a correction factor (following Ferguson, 1986) may need to be applied. We will provide the value of the unbiased estimator of the variance ( $s$  in Ferguson, 1986) in the text. See also our response to comment 4 of Anonymous Reviewer 2.

The following text was added: “Because these corrections are being carried out in log10 space statistical bias (underestimation) is possible (Ferguson, 1986). Following Ferguson (1986) we calculate the unbiased estimate of the variance ( $s$ ) as 0.0686 which gives a correction factor (calculated as  $\exp(2.65s^2)$  of 1.0126 may be applied uniformly as a correction factor (multiplier) to the bias corrected values to adjust for possible bias due to the log10 space regression we have implemented.”

(4) Other comments: L152-153: Perhaps, a more appropriate way of describing  $NSE \leq 0$  would be that the "model is no better than using the mean value of the observed data as a predictor" (e.g. Gupta et al. 2009, Journal of Hydrology). This conveys a somewhat different meaning than how it is presently written, i.e. “values less than zero indicate that the mean value of the observed data is a better predictor than the hydrologic model”

Agreed, as written it was inappropriate and convoluted. We will rewrite this sentence to read: “NSE values are useful (compared to VE) in that values less than zero indicate that the model is no better than using the mean value of the observed data as a predictor”

L196-197: Or alternatively, that model-data agreement improves with catchment size

Excellent point, that is likely due to training on GRUN’s large catchments elsewhere. Thus, we will rewrite this sentence to say both: “This suggests two possibilities: first, that particularly for small catchments which may have steeper average slope, GRUN underpredicts monthly runoff values associated with the wet season; and second model-data agreement improves with catchment size. Bias is likely inherited from the high uncertainty in monsoonal precipitation rates inputted into GRUN.” This point is also discussed later in the paper.

L226-227: “These catchments experience distinct wet and dry seasons in the northwest Philippines.” Can the authors comment on the implication of this sentence for catchments (outside the Philippines) with distinct wet and dry seasons vis-à-vis the physical significance (e.g. of rainfall-runoff transfer functions) that the VE criterion represents?

This is a question of baseflow not being properly represented well. For large basins where the time of concentration ( $T_c$ ) of any given flood event will be much longer, Flood peaks may be wider and subdued because of abstractions and infiltration into the shallow aquifer. This phenomenon is significantly less apparent in smaller basins where peak flows are expected to be higher because of less infiltration and narrower peaks. We will add this explanation to the revised manuscript.

These two sentences were added: “The physical significance of these observations are that for large basins the time of concentration of any given flood event will be much longer, thus flood peaks will be wider and subdued due to abstractions and infiltration into the shallow aquifers. This phenomenon is likely less apparent in smaller basins where peak flows are expected to be higher because of less infiltration.”

L230-231: Please qualify/rewrite because while Criss and Winston (2008) underlined that NSE tends to put more weight on large flows, they did not particularly discuss or say anything regarding NSE-log10.

Thank you for pointing this out. We will instead say: “Since NSE puts more weight on large flow (Criss and Winston, 2008), it is not surprising that our NSE-log10 scores are in most cases significantly more skillful among our catchments because extremely wet months are weighted less than using the raw runoff values compared to raw NSE scores.”

L240-243: Please consider rewriting this long sentence for clarity. Also, this sentence refers particularly to GRUN (published in 2019) yet it cites two papers that predates GRUN. Please qualify for congruence.

Agreed. We will change to two sentences and change the parenthetical citation to: “As such, we suggest that GRUN, is a useful new tool for studying trends, seasonality and average runoff from tropical catchments (such as in previous work: e.g., Merz et al., 2011; Wanders and Wada, 2015) in the Philippines. However, we qualify this finding by noting that GRUN is not suitable for extreme value analyses associated with major tropical storms during the wet seasons unless suitable bias corrections (see next section) can be effectively carried out.”

Technical corrections

L61: “(ref)”. Reference placeholder

Thank you, this was a mistake on our part, we will cite: Hagemann et al. (2011, *J. of Hydrometeorology*), Davie et al. (2013, *Earth System Dynamics*) and Winsemius et al. (2016, *Nature Climate Change*) as examples in this parenthetical. Other examples welcome.

L99: URL is not working. Please check

Thank you for checking, unfortunately this URL has gone defunct since submission, we are contacting the DPWH who provided a new URL; however as of August 2020 this link is not yet working. All data used in this analysis is provided in our supplemental files. We will attempt to rectify this issue with the reporting agency during the proof stage.

L126: “that”.

Will be removed.

Typo L206: “were”. Typo

Will be changed to “where”

## Response to Anonymous Referee #2

Ibarra et al. compare the monthly streamflow data from a new global database (GRUN) with gauge data for multiple catchments in the Philippines. The catchments were not included in the development of the global database and are smaller than those used to train the machine-learning algorithm to create the global database. The work is very interesting because it highlights issues with the use of such global databases for smaller catchments or to answer local questions.

I thank and congratulate the authors for putting together a very valuable database of streamflow for catchments in the Philippines. This is highly useful because so much of our collective research efforts and knowledge focuses on catchments in temperate climates.

Nonetheless, I do have several comments (see below). These mainly focus on the lack of a comparison with a lower (and upper) benchmark, the use of overall statistics and the effects of pooling of data with different record lengths and variability on these overall statistics, the over-use of log-log transformation and scales, the overall message, and the appropriateness of the technical note category.

Thank you for the encouraging, constructive and thorough review. We respond below to each of these points in detail.

I include other detailed comments and suggestions to strengthen the manuscript in the attached pdf. This pdf also has some editorial suggestions. Note that these are just suggestions to make the text more to the point or clearer - I don't think that all of them need to be addressed or implemented.

We thank the reviewer for edits in the attached pdf, all of these editorial suggestions will be taken/addressed in our revision.

1. To assess the value and skill of the GRUN database, the results need to be compared to a lower benchmark. How well (or poor) does a very simple estimate reflect the observed monthly streamflow and how much better are the GRUN estimates than this rough estimate or lower benchmark? This lower benchmark could be based on the multiplication of the average runoff ratio for the region and month by the local rainfall data or the average streamflow from all catchments in a region. Similarly, we can not expect a perfect  $r^2$  or NSE score because the observations are uncertain. Although the uncertainties are mentioned for the data in one database, they are not shown or discussed anywhere. As a result it remains difficult to interpret the results, i.e., should we consider

these  $r^2$ , VE, NSE or RMSE values as poor or does the GRUN data have some reasonable predictive power that is much better than a regional average value?

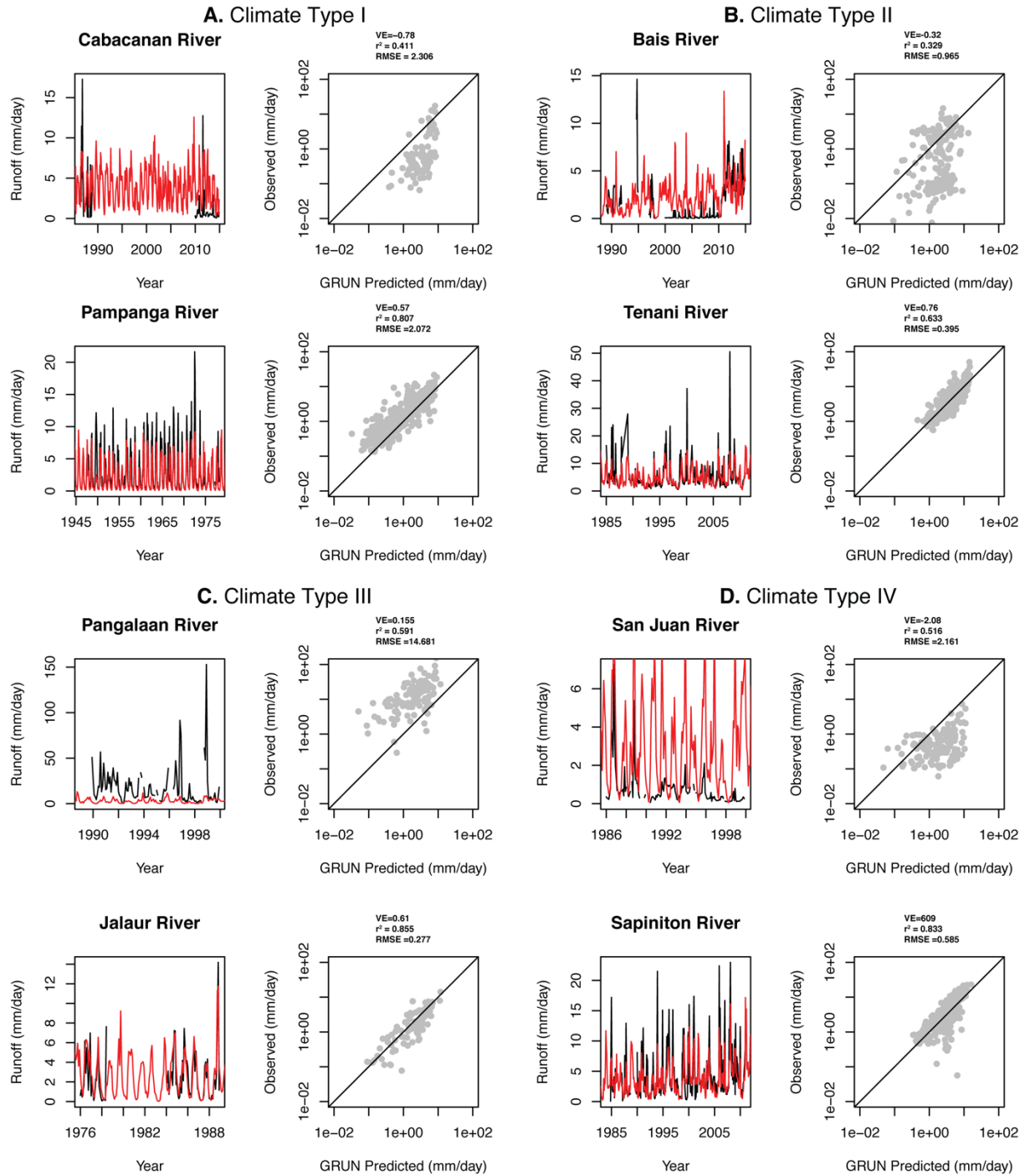
Regarding the lower benchmark, we view the NSE as a good estimator of a lower benchmark. Given that the NSE value before bias correction is quite poor (country-wide) with a value of 0.091 (and a log10 value of 0.453), and an improvement in this score following bias correction, we view this as significantly better than the lower benchmark of just using the mean value of the observed data. Unfortunately, we have a limited number of paired rain gauge and catchment discharge data to perform a runoff ratio calculation as suggested by the reviewer. However, the use of satellite-based rainfall estimates calibrated with limited rain gauge data coupled to GRUN discharge data is the subject of another paper our group is currently preparing.

Regarding the metrics used: as a point of comparison the  $r^2$  and NSE values that are calculated and reported in our Table 2, are lower than again the large GRDC river basins compared to in the original GRUN paper (Ghiggi et al., 2019; see their Table 2). This is an important point of our paper, even once bias corrected using the compiled datasets the NSE value is lower than that of the global comparison. This is important as it demonstrates the need to incorporate smaller tropical catchments such as those presented here in such global runoff datasets, particularly with respect to correctly predicting high-flow during wet seasons. In addressing other comments from this reviewer and the other reviewer we have now provided this context and described the issues with GRUN data (as currently presented). We are now working with the GRUN lead author (G. Ghiggi) to provide enhanced data coverage from these data (and earlier) to improve the next iteration of GRUN, and hope that the qualifications we have added (pasted in below responses) in sections 2.1 and the conclusions address these concerns.

2. Although the time series of observed and GRUN based streamflow are given in the supplementary material, none of these time series are shown in the main document. I highly recommend showing these plots for the best and the poorest site in each climate zone in the main document. This will give the reader a much better feeling of the skill of the GRUN data and helps with the interpretation of the VE,  $r^2$  and RMSE values that are given in the text.

To address this comment we will add a new Figure (between figures 1 and 2 currently) showing 8x examples from the existing supplementary figures (time series and cross plots), listing the VE,  $r^2$  and RMSE values for that catchment.



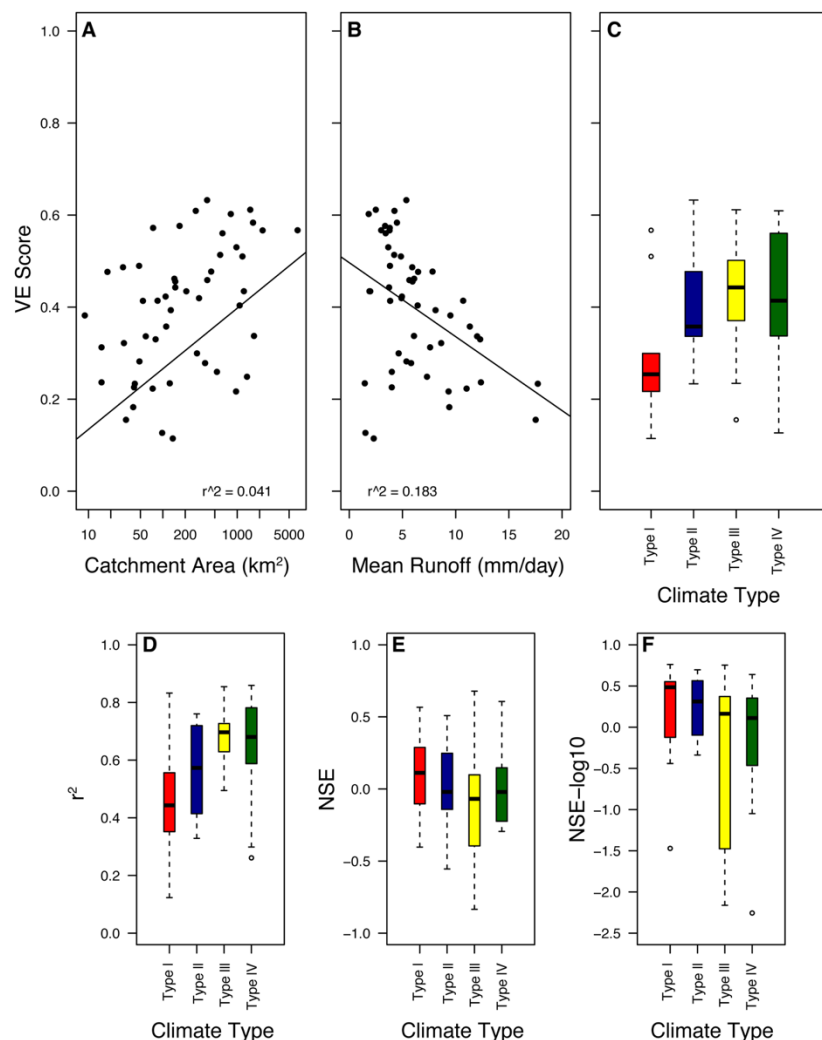


3. A large part of the analyses is based on pooled data (e.g. Figure 3) but the record length is very different for the different catchments and the variability in discharge is also different. This likely influences these pooled results. I would rather (also) see boxplots that show the  $r^2$ , NSE and NSE-log values for the individual catchments as is done for VE in Figure 4. This will also give the reader a much better feeling of how different these results are for the different catchments. I therefore suggest to add these plots to the manuscript and to add these ranges to Table 2 as well.



We agree to the comment thus we placed importance in including individual metrics so readers can see and reinterpret biases in the dataset as what reviewer 2 has done. We will include similar box plots of  $r^2$ , NSE and NSE-log values for the individual catchments as in Figure 4 for VE and add additional lines with the median and IQR to Table 2. We append a draft of this figure to this reply, this will be added to Figure 4.

Note that for clarity we kept the VE score in the top set of panels and then put box plots of  $r^2$ , NSE and NSE-log values in a second set of panels below. See revision for new figure (now Figure 5 due to addition of figures).



- Almost all comparisons of the observed and GRUN-based discharge are shown in log-log space. This is informative for some analysis and allows one to see the data but at the same time almost any comparison looks OK in log-log space, even when the data don't really match. I therefore suggest to not use log-log axes where it is not entirely necessary.

For example figure 1a could be split in 3 sub-panels (max, median,min) and then show the data on a linear scale. Furthermore, I wonder whether the bias correction in log-space leads to large errors when the corrected values are transferred again. This isn't shown nor discussed in the manuscript.

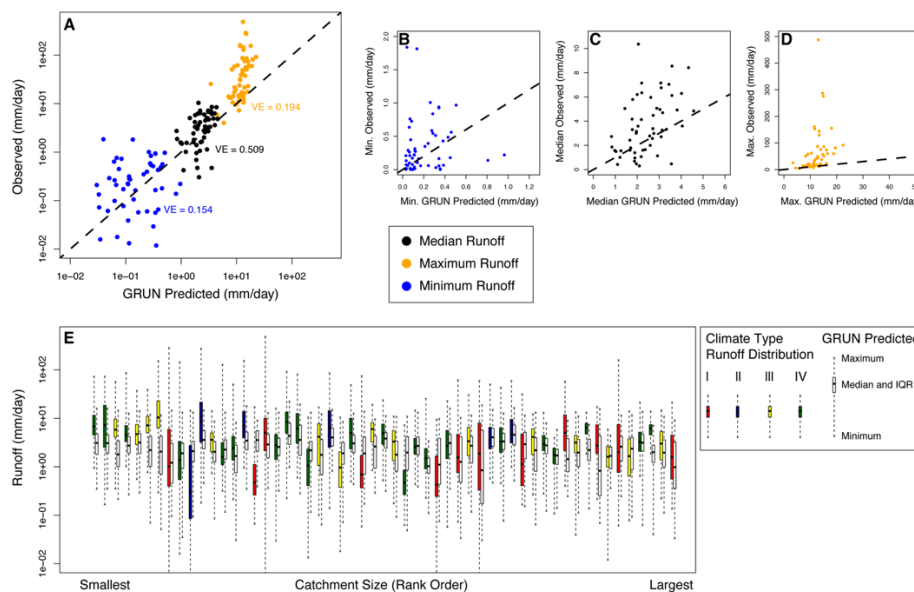
This is a good comment that was also brought up by the first reviewer. To rectify we will:

- Split Figure 2a (now 3a) into 3 sub-panels in linear scale for the max, median and minimum runoff.
- Bias correction via log-space does introduce errors, particularly for large runoff values with ratio of the true to predicted values scaling approximately as  $\exp(2.65s^2)$  where  $s$  is the unbiased estimator of the variance (see Ferguson (1986) referred to us by the first reviewer). To account for this we will both perform a bootstrap analysis (shown above) and then mention in the text that it is suggested to also apply this scaling and provide the value for the unbiased estimate of the variance ( $s$ ).

Ferguson, R. I. River Loads Underestimated by Rating Curves. Water Resour. Res. 22, 74–76 (1986).

See response above to the first reviewer for more details and the text that was added: “Because these corrections are being carried out in log10 space statistical bias (underestimation) is possible (Ferguson, 1986). Following Ferguson (1986) we calculate the unbiased estimate of the variance ( $s$ ) as 0.0686 which gives a correction factor (calculated as  $\exp(2.65s^2)$  of 1.0126 may be applied uniformly as a correction factor (multiplier) to the bias corrected values to adjust for possible bias due to the log10 space regression we have implemented.”

The new figure (now Figure 3)



5. I think that the overall conclusion of the manuscript is too optimistic (although this admittedly depends on the comparison with the lower benchmark (see comment 1)). I agree that the bias correction helps and leads to a significant improvement but the GRUN based estimates of streamflow are very poor (particularly when one looks at the time series that are given in the supplementary material). The abstract and conclusion could highlight the danger in using these types of products for local streamflow prediction, model calibration, etc more clearly. Now it seems to be overly optimistic on how these data can be used for a range of local studies or to answer local questions.C3

We agree that our overall conclusions are overly optimistic, we will add the following to the conclusions: “Global databases such as GRUN are applicable for aggregated stream discharge estimates and to investigate general trends in the hydrologic characteristics of a region. The recommended bias correction presented here will likely improve such estimates and analysis for the Philippines. While GRUN was never intended to be used for estimating single catchment discharge its applicability for such purposes can be extended provided that proper statistical comparison of modelled versus actual gauged data are initially performed. We thus propose that the utilization of the GRUN dataset can be extended to other ungauged tropical regions with smaller catchments upon applying a similar correction as described in this study.”

6. The manuscript was submitted as a technical note but the manuscript doesn't fully fit the description of a technical note as it doesn't describe a new method or technique. It should thus be a regular research paper, which would allow for more comparisons of the datasets (as described above) and additional figures (as described above).

We have consulted with the HESS editor and do feel that this does fit the description of a technical note because we have compiled a new runoff dataset from a previously unrepresented country for the community to use and provided a tool via our bias correction equation for correcting a global observation-based gridded runoff dataset for use based on this data. Based on our email exchange with the HESS editor (Wouter Buytaert) we have not changed the classification of this contribution.

7. Some more background information on the gauging station data used in GRUN would be helpful, e.g. what percentage of stations were located in the tropics? And do the papers that describe the GRUN database make claims about smaller catchments or tropical catchments?

We agree with the reviewer that these are important details regarding GRUN that we overlooked mentioning, we will add the following details to the revised manuscript:

- GRUN is highly biased to the northern hemisphere mid-latitudes. We will refer the reader to the original publication (Ghiggi et al., 2019) with respect to details mentioning the underrepresentation of small tropical catchments (not just due to the criteria of  $>10,000$  km<sup>2</sup> but also the lack of records in general).

- Additionally, we will mention that the GRUN publication does discuss that because of how they trained the dataset the uncertainty is greatest the tropics (see bottom left of page 1164 and Figure 8b of Ghiggi et al., 2019) in comparison to other datasets, and mention southeast Asia showing an increase in runoff rates over the period of analysis despite a paucity of records (see top right of page 1666).
- Further, a strong correlation of runoff with ENSO is shown in figure 11 of Ghiggi et al. (2019) in southeast Asia, though the resolution of the map makes it difficult to ascertain significance and regional variability in the Philippines.

In consultation with Gionata Ghiggi the following was added to the methods: “Additionally, due to the training criteria GRUN’s calibration is biased towards an overrepresentation of the northern hemisphere mid-latitudes relative to the tropics with few sites available in Africa and southeast Asia. Ghiggi et al. (2019) discuss that because of the dataset training techniques uncertainty scales with the magnitude of runoff rates and is likely to have high prediction uncertainty in regions with less dense runoff observations such as in tropical southeast Asia. Further, they show in southeast Asia an increase in runoff rates and a strong correlation of runoff with ENSO over the period of analysis (1902 to 2014). We refer the reader to Ghiggi et al. (2019) for more information but note that because of the catchment size filtering criteria, none of the GISM and GRDC data from the Philippines was used (Personal Communications, G. Ghiggi, 2019). As such, we view our analysis as a completely independent test of the GRUN runoff reconstruction using small tropical catchments.”

We thank the reviewer for these questions about GRUN, while not our group’s product it will certainly strengthen our manuscript.

# Technical Note: Evaluation and bias correction of an observations-based global runoff dataset using historical streamflow observations from small tropical catchments in the Philippines

Daniel E. Ibarra<sup>1,2</sup>, Carlos Primo C. David<sup>3</sup>, Pamela Louise M. Tolentino<sup>3</sup>

<sup>1</sup> Department of Earth and Planetary Science, University of California, Berkeley, California 94720 USA

<sup>2</sup> Institute at Brown for Environment and Society and the Department of Earth, Environmental and Planetary Science, Brown University, Providence, Rhode Island 02912 USA

<sup>3</sup> National Institute of Geological Sciences, University of the Philippines, Diliman, Quezon City, Philippines 1101

Correspondence to: Daniel E. Ibarra ([dibarra@berkeley.edu](mailto:dibarra@berkeley.edu)) and Carlos Primo C. David ([cpdavid@nigs.upd.edu.ph](mailto:cpdavid@nigs.upd.edu.ph))

## Abstract

The predictability of freshwater availability is one of the most important issues facing the world's population. Even in relatively wet tropical regions, seasonal fluctuations in the water cycle complicate the consistent and reliable supply of water to urban, industrial and agricultural demands. Importantly, historic streamflow monitoring datasets are crucial in assessing our ability to model and subsequently plan for future hydrologic changes. In this technical note we evaluate a new global product of monthly runoff (GRUN; Ghiggi et al., 2019) using small tropical catchments in the Philippines. This observations-based monthly runoff product is evaluated using archived monthly streamflow data presented in this study from 55 catchments with at least 10 years of data, extending back to 1946 in some cases. Since GRUN didn't use discharge data in the Philippines to train/calibrate their models, data presented in this study provide the opportunity to evaluate independently such product. We demonstrate across all observations significant but weak correlation ( $r^2 = 0.372$ ) and skilful prediction (Volumetric Efficiency = 0.363 and log(Nash-Sutcliffe Efficiency) = 0.453) between the GRUN predicted values and observed river discharge. At a regional scale we demonstrate that GRUN performs best among catchments located in Climate Types III (no pronounced maximum rainfall with short dry season) and IV (evenly distributed rainfall, no dry season). We also find a weak negative correlation between volumetric efficiency and catchment area, and a positive correlation between volumetric efficiency and mean observed runoff. Further, analysis of individual rivers demonstrates systematic biases (over and under) in baseflow during the dry season, and under-prediction of peak flow during some wet months among most catchments. These results demonstrate the potential utility of GRUN and future data products of this nature with due consideration and correction of systematic biases at the individual basin level to: 1) assess trends in regional

Deleted: \_v1

Deleted: These catchments are completely independent of the

Deleted: gridded product as no catchments

Deleted: were of sufficient size

Deleted: fulfil the original filtering criteria and databases of these data were either not digitized or difficult to compile. Using monthly runoff observations from catchments with more than 10 years of...

Deleted: between 1946 and 2014, we

Deleted: the observations

40 scale runoff over the past century, 2) validate hydrologic models for un-monitored catchments in the Philippines, and 3)  
assess the impact of hydrometeorological phenomenon to seasonal water supply in this wet but drought prone archipelago.  
Finally, to correct for underprediction during wet months we perform a log-transform bias correction which greatly improves  
the nationwide Root Mean Square Error between GRUN and the observations by an order of magnitude (2.648 vs. 0.292  
mm/day). This technical note demonstrates the importance of performing such corrections when accounting for the  
45 proportional contribution of catchments from smaller catchments in tropical land such as the Philippines to global tabulations  
of discharge.

## 1 Introduction

The global water crisis is considered as one of the three biggest global issues that we need to contend with, affecting  
an estimated two-thirds of the world's population (Kummu et al., 2016; WEF, 2018). Among the sources of freshwater, the  
50 most important compartment in terms of utility is surface water flow, which is the primary resource for irrigation, industrial  
use and for bulk domestic water supply for many large cities. Along with the purpose of flood mitigation during extreme  
weather events, monitoring streamflow is a vital activity that many nations conduct with various levels of coverage. Long  
term streamflow datasets prove useful in resource management and infrastructure planning (e.g., Evaristo and McDonnell,  
2019). Such data is even more critical in areas that rely on run-of-the-river supply and do not utilize storage structures such  
55 as dams and impoundments. Further, a robust, long term dataset is crucial in the face of increased variability in stream  
discharge due to land use change, increased occurrence of mesoscale disturbances and climate change (e.g., Abon et al.,  
2016; David, et al., 2017; Kumar et al., 2018).

There is a disparity in the availability of long-term gauged rivers datasets between continental areas and smaller  
island nations. This in the face of the latter having an invariably more dynamic hydrometeorologic system owing to the size  
60 of their catchment and proximity to the ocean (e.g., Abon et al., 2011; Paronda et al., 2019). Furthermore, in the case of  
tropical island nations these are where the impact of climate change in the hydrologic cycle could be observed the most  
(Nurse et al., 2014). Thus, the Philippines offers a unique example where manual stream gauging programs have started in  
1904 and, while spotty at times, have continued on to today. In this work we analyse data since 1946. This island nation on  
the western side of the Pacific Ocean shows a very dynamic hydrologic system as affected by tropical cyclones, seasonal  
65 monsoon rains, sub-decadal cycles such as the El Nino Southern Oscillation (ENSO) and overlaid on top of all these are the  
hydrologic changes caused by climate change (Abon et al., 2016; David, et al., 2017; Kumar et al., 2018).

In the absence of long-term streamflow datasets, several researchers have compiled datasets worldwide which are  
used to extrapolate streamflow in non-gauged areas (Maybeck et al., 2013; Gudmundsson et al., 2018; Do et al., 2018;  
Alfieri et al., 2020; Harrigan et al., 2020). Several global hydrological models have also been created to project variations in  
70 streamflow and extend present-day measurements to the future (Hagemann et al., 2011; Davie et al., 2013; Winsemius et al.,  
2016). The latest contribution to modelled global runoff products is the Global Runoff Reconstruction (GRUN) (Ghiggi et

Deleted: ).

Deleted: ref).

Deleted: streamflow product

75 al., 2019). GRUN is a global gridded reconstruction of monthly runoff [for the period 1902-2014](#) at 0.5 degree (~50km by 50km) spatial resolution. It uses global streamflow data from [7,264](#) river basins that to train a machine learning algorithm [which learn the runoff generation processes from](#) precipitation and temperature data.

This technical note evaluates the accuracy of the GRUN dataset (GRUN\_v1) as applied to the hydrodynamically-active smaller river basins in the Philippines. Additionally, it explores the possible hydrologic parameters that may need to  
80 be considered and/or optimized such that global datasets may be able to model runoff in smaller, ungauged basins more accurately.

- Deleted: 718 large
- Deleted: is used
- Deleted: that would be able to take
- Deleted: to predict monthly global runoff for the period 1902-2014.

## 2 Dataset and Methods

### 2.1 GRUN observations-based global gridded (0.5°x0.5°) runoff dataset

85 GRUN is a recently published global reconstruction of monthly runoff time series spanning 1902 to 2014 created using a machine learning algorithm based on training temperature and precipitation fields from the Global Soil Wetness Project Phase 3 (GSWP3; Kim et al., 2017; <http://hydro.iis.u-tokyo.ac.jp/GSWP3/index.html>) using the Global Streamflow Indices and Metadata Archive (GISM) (Ghiggi et al., 2019). In this contribution we [analyse GRUN\\_v1](#) ([https://figshare.com/articles/GRUN\\_Global\\_Runoff\\_Reconstruction/9228176](https://figshare.com/articles/GRUN_Global_Runoff_Reconstruction/9228176); accessed September 9<sup>th</sup> 2019) which was  
90 trained on [a selection of small catchments with area between 10 and 2500 km<sup>2</sup>](#) [GISIM](#) (Do et al., 2018; Gudmundsson et al., 2018) and validated using [379](#) large (>50,000 km<sup>2</sup>) monthly river discharge datasets from the Global Runoff Data Centre (GRDC) Reference Dataset ([https://www.bafg.de/GRDC/EN/04\\_speltdbss/43\\_GrFN/refDataset\\_node.html](https://www.bafg.de/GRDC/EN/04_speltdbss/43_GrFN/refDataset_node.html)). [Additionally, due to the training criteria GRUN's calibration is biased towards an overrepresentation of the northern hemisphere mid-latitudes relative to the tropics with few sites available in Africa and southeast Asia. Ghiggi et al. \(2019\) discuss that because of the dataset training techniques uncertainty scales with the magnitude of runoff rates and is likely to have high prediction uncertainty in regions with less dense runoff observations such as in tropical southeast Asia. Further, they show in southeast Asia an increase in runoff rates and a strong correlation of runoff with ENSO over the period of analysis \(1902 to 2014\).](#) We  
95 refer the reader to Ghiggi et al. (2019) for more information but note that because of the catchment size filtering criteria, none of the GISM and GRDC data from the Philippines was used (*Personal Communications*, G. Ghiggi, 2019). As such, we  
100 view our analysis as a completely independent test of the GRUN runoff reconstruction using small tropical catchments.

- Deleted: analyze GRUN\_v1
- Deleted: <https://doi.org/10.3929/ethz-b-000324386>
- Deleted: large (> 10,000
- Deleted: ) stations from GISM
- Deleted: 718

### 2.2 Historical streamflow observations

In this contribution we analyse monthly observations of discharge from [55](#) manually observed streamflow stations from three Philippine datasets. The observations span 1946 to 2016, although only data through 2014 are utilized due to the  
105 time period included in GRUN.

- Deleted: 74
- Deleted: \_v1.
- Formatted: Indent: First line: 0"



2.2.1 Bureau of Research and Standards (BRS) Dataset

The historical discharge data was originally acquired from the Bureau of Research Standards (BRS) under Department of Public Works and Highways (DPWH). The records keeping was transferred to the Bureau of Design, also under DPWH, which continues to record gage data from some rivers up to this date. The degree of accuracy of records were categorized as “excellent”, “good”, “fair”, or “poor” using the following convention: “Excellent” means about 95% of daily discharges are within ±5% difference of the actual gauge height vs height computed from the rating curve; “Good” is within ±10%; and “Fair” is within ±15%; while “Poor” means daily discharges are below the 15% “Fair” accuracy. This is the only basis of accuracy for this set of data. A majority of the reprocessed BRS data used in this analysis come from Tolentino et al. (2016), however, some of the datasets were subsequently updated using [data available from the Department of Public Works and Highways](https://apps.dpw.gov.ph/streams_public/station_public.aspx).

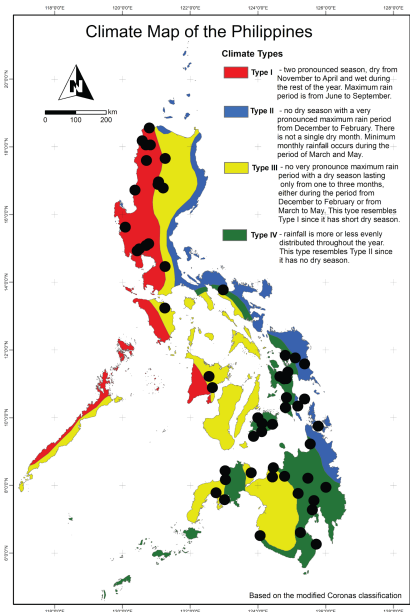
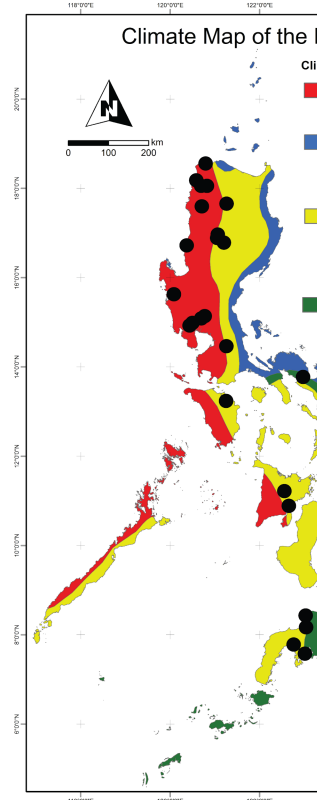


Figure 1: Map of Philippines with location of streamflow stations used in this analysis overlaid on climatic type (as in Tolentino et al., 2016; Kintanar, 1984; Jose and Cruz, 1999). Note that no publicly available long-term stations are available for Palawan.

Deleted: the online archive for this study ([https://apps.dpw.gov.ph/streams\\_public/station\\_public.aspx](https://apps.dpw.gov.ph/streams_public/station_public.aspx))...



Deleted:

140 **2.2.2 Global Runoff Data Centre (GRDC) Reference Dataset**

Ten catchments from the GRDC Reference Dataset were compiled (https://www.bafg.de/GRDC/EN/04\_spcldtbss/43\_GrFN/refDataset\_node.html; requested July 2019) and analysed in this contribution. Over 45 sites from the Philippines are available in the GRDC data; however, almost all do not fulfil our criteria of having over 10 years of coverage. Four of these catchments also duplicated or extended the BRS datasets, and one extends a GSIM dataset (see below). Notably four of the time series available from GRDC are available back to the 1940s (Table 1).

150 **2.2.3 Global Runoff Data Centre (GSIM) Reference Dataset**

Only two time series of the available five GSIM time series (Gudmundsson et al., 2016, Do et al., 2018) cover more than 10 years.

**2.3 Criteria for Inclusion of Datasets**

All catchment areas were verified using the digital elevation model from the 2013 Interferometric Synthetic Aperture Radar (IFSAR) data. All hydrologic datasets were normalized to runoff (mm/yr), sometimes also notated as ‘specific discharge’ in the literature. We only considered streamflow stations where the published and verified areas agreed and coverage spanned 10 or more years. The location of all streamflow stations is shown on Figure 1 and listed in Table 1. Catchment areas span 4 orders of magnitude (8.93 to 6487 km<sup>2</sup>) and cover the majority of the Philippines excluding Palawan (see Figure 1). The location of catchments was paired to GRUN grid cells (0.5° by 0.5°) for the analysis. Instead of computing the weighted area runoff over the catchment, we employed nearest neighbour interpolation between the catchment outlet location and the GRUN gridded product (0.5° by 0.5° resolution). All but one catchment is smaller than the area of the GRUN grid cells (~2,500 km<sup>2</sup>), thus, we view this pairing as sufficient for validation purposes. This assumption was tested by interpolating the GRUN grid to the gauging location as well as the watershed centroids, and no significant difference in correlation to the observations that were observed.

**2.4 Statistical Performance Metrics**

165 To assess the performance of GRUN we use a suite of metrics commonly used to assess model performance in hydrologic studies. These metrics are calculated for each individual catchment (n=55) and for each climate type (n=4; see below) shown in Figure 1.

170 Firstly, we use the commonly used square of the Pearson correlation coefficient ( $r^2$ ). This metric for bivariate correlation measures the linear correlation between two variables. In this case the predicted monthly values from GRUN versus the observed monthly values from the streamflow datasets. It varies from 0 (no linear correlation) to 1 (perfect correlation). The use of  $r^2$  does not account for systematic over- or under-prediction in runoff because it only accounts for

**Deleted: 1**

**Deleted:** Maybeck et al., 2013,

**Deleted:** Further, none of these data were included in the original training of GRUN due to catchment size constraints.

**Deleted:** Given that all

**Deleted:** 3,000

**Deleted:** )

correlation among the observed and predicted values (see Krause et al. (2005) for further discussion of the use of  $r^2$  in hydrological model assessment).

Secondly, the primary metric used here is Volumetric Efficiency (VE) (Criss and Winston, 2008), utilized previously by Tolentino et al. (2016) on a subset of the BRS catchments analysed here. VE is defined as:

$$VE = 1 - \frac{\sum Q_P - Q_O}{\sum Q_O} \tag{1}$$

where P is the modelled/predicted values and O is the observed runoff values. A value of 1 indicates a perfect score. Because we are interested in the performance of GRUN over the period of each streamflow record, unlike Tolentino et al. (2016), we calculate VE using all paired monthly observed and simulated values rather than the monthly medians. This results in lower VE scores than those previously reported by Tolentino et al. (2016) compared to hydrologic models.

Further, we use both the linear and logged Nash-Sutcliffe efficiency (NSE) parameter. Proposed by Nash and Sutcliffe (1970) it is defined as:

$$NSE = 1 - \frac{\sum (Q_P - Q_O)^2}{\sum (Q_O - \text{mean}(Q_O))^2} \tag{2}$$

The range of NSE can be between  $-\infty$  and 1 (perfect fit). NSE values are useful (compared to VE) in that values less than zero indicate that the model is no better than using the mean value of the observed data as a predictor. NSE is also calculated using logarithmic values prior to calculation to reduce the influence of peak flow and increase the influence of low flow values (see further discussion in Krause et al., 2005).

Finally, to evaluate a possible strategy for performing a bias correction of the GRUN simulated values at a countrywide scale we use the Root Mean Square Error (RMSE) in units of runoff (i.e., mm/day). The RMSE is applied to the raw GRUN simulated values and the observation based bias corrected GRUN values at the country, climate type (see below) and individual catchment level.

### 2.5 Climate Types

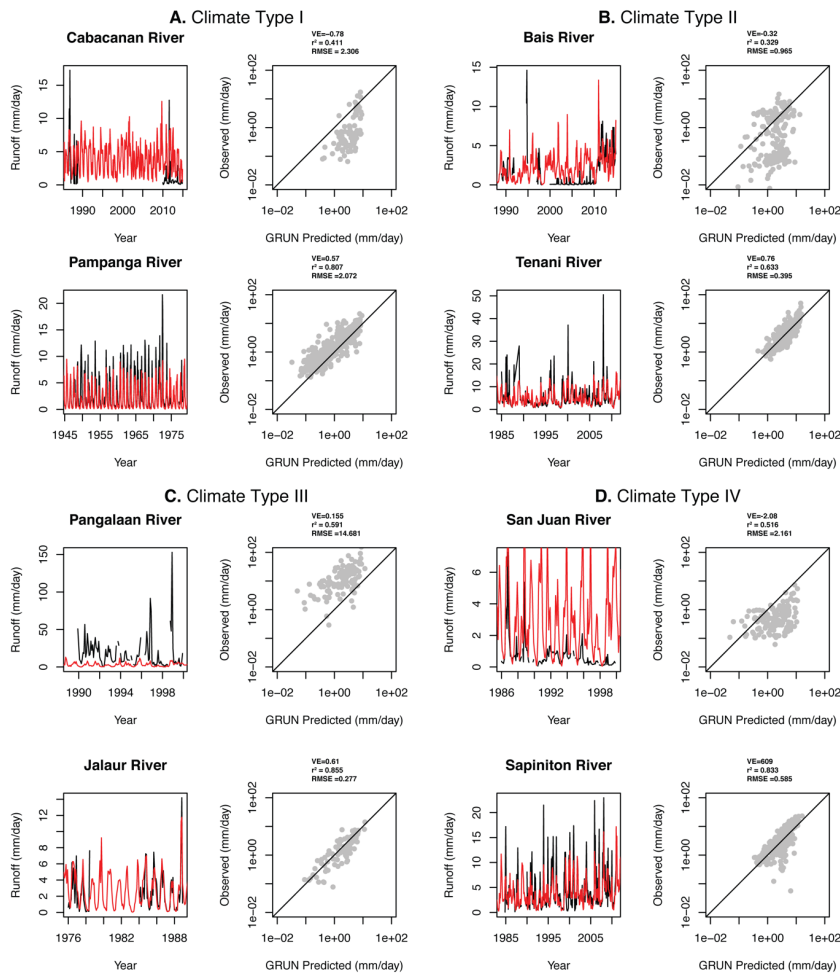
The Philippines has four Climate Types (see also Abon et al., 2016; Tolentino et al., 2017): Type I Climate on the western seaboard of the Philippines is characterized by distinct wet (May to October) and dry (November to April) seasons; Type II Climate on the eastern seaboard has no distinct dry period with maximum rains occurring from November to February; Type III inland climate experiences less annual rainfall with a short dry season (December to May) and a less pronounced wet season (June to November); and, Type IV southeast inland climate experiencing depressed rainfall and characterized by an evenly distributed rainfall pattern throughout the year. Further description is provided in Figure 1.

Deleted: is

Deleted: better

Deleted: than the hydrologic model

Deleted:



Deleted:

**Figure 2.** Example timeseries of GRUN predicted (red lines) and observed (black lines) runoff values, and cross-plots (log-scale) with VE,  $r^2$  and RSME values for the worst (top) and best (bottom) performing river basins within Climate Types I, II, III and IV (panels A-D, respectively).

### 3 Results and Discussion

Statistical comparisons described above between individual catchments and the GRUN dataset are shown in the supplemental figures and tabulated in Table 1. Shown [as in examples in Figure 2 and](#) also in the supplemental figures (Figure A1) are time series comparisons between the GRUN runoff values and runoff (area normalized discharge). Statistics performance metrics across all data as well as by climate types I to IV are also listed in Table 1. Given the emphasis on a country scale evaluation of GRUN we primarily focus below on results in aggregate grouped by climate type or among all catchments.

Across all observations GRUN has a r-squared of 0.372 and a VE of 0.363 (Table 2). Using log(runoff) values (following Criss and Winston, 2008) this improves to an r-squared of 0.546 and a volumetric efficiency (VE) of 0.733, suggesting reasonable utility in the GRUN product at a country scale for the Philippines, despite no training data from the Philippines being used in the creation of GRUN. The raw RMSE across the dataset is 2.648 mm/day (Table 2). In the following we break down the comparison between the streamflow observations and GRUN by first comparing runoff distributions and extreme values at the individual basin level, then aggregating our results by Climate Type, and finally look at several correlations of VE to watershed characteristics.

#### 3.1 Comparison of runoff distributions

Average runoff values among all catchments compared to GRUN show reasonably good predicted values. In Figure 3A median values of runoff (black dots) given a volumetric efficiency (VE) metric of 0.509 across all catchments, the average (mean) difference between median observations and simulated values is +16%. Looking at extreme monthly values (maximum and minimum) over the months of observation demonstrates significant underprediction in wettest conditions (orange dots in Figure 3A and 3D) with almost all catchments' maximum observations falling above the 1:1 line and a lower VE score of maximum values of 0.194. For baseflow conditions spread around the 1:1 line for minimum values is more evenly distributed, however the VE score of 0.154 is similarly low due to a greater spread than the median values.

This suggests two possibilities: first, that particularly for small catchments which may have steeper average slope, GRUN underpredicts monthly runoff values associated with the wet season; and second model-data agreement improves with catchment size. Bias is likely inherited from the high uncertainty in monsoonal precipitation rates inputted into GRUN. However, we do find a significant correlation (at  $p < 0.01$ ,  $r^2 = 0.391$ ) between log values of maximum runoff difference (observed minus predicted) versus catchment area (not shown), with a negative relationship.

In general, the median and interquartile ranges (IQR, 25% to 75%) shown in Figure 2E overlap between GRUN and the observations. For five catchments of large ( $n=2$ ) and relatively small sizes ( $n=3$ ) the IQR of the observations does not overlap with the GRUN runoff IQR. The three small catchments are Climate Type III (yellow) and the two large catchments are Climate Type IV. In two catchments of moderate size the GRUN IQR is greater than the observed IQR runoff range.

**Deleted:** (Personal Communications, G. Ghiggi, 2019).

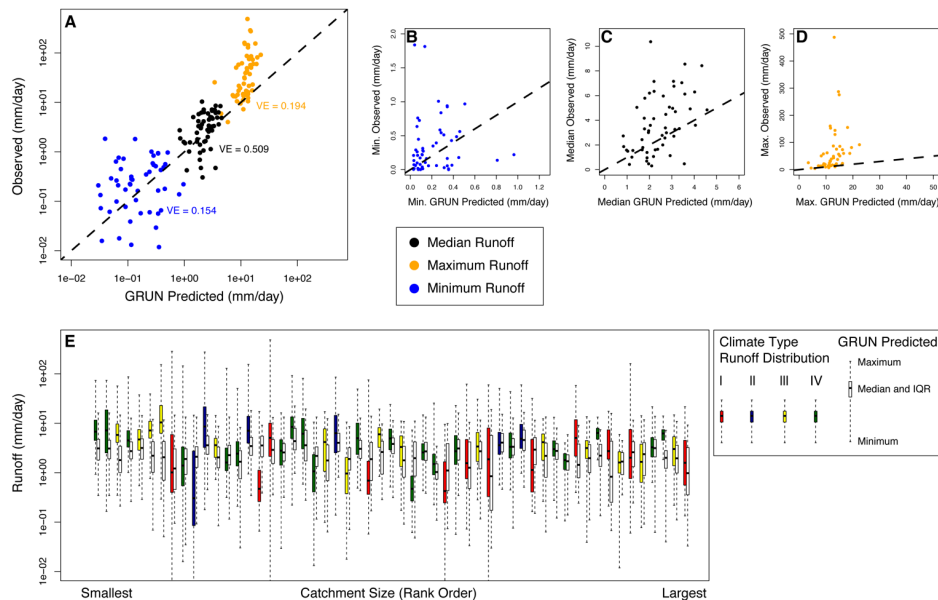
**Deleted:** 2A

**Deleted:** 2A

**Deleted:** → Comparison of runoff distributions ranked by catchment size in Figure 2B demonstrates that maximum runoff values appear to be most underpredicted by GRUN in the smallest catchments. →

**Deleted:** This suggests that particularly for small catchments, which may have steeper average slope, GRUN underpredicts monthly runoff values associated with the wet season. ¶

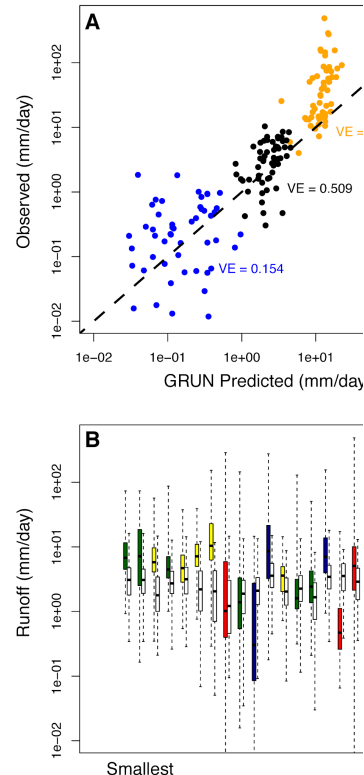
**Deleted:** 2B



**Figure 3.** Comparison of runoff ranges and distributions. (A) Comparison of median and extreme (maximum and minimum) monthly values between the observations and GRUN, in log space. (B-D) As in (A) for the minimum, median and maximum monthly values, respectively, in linear space. (E) Distribution of runoff between observations (coloured) and GRUN (white) using box and whisker plots. Plots show the median, interquartile range and maximum/minimum values. The GRUN distributions only include months where observations are present.

### 3.2 Comparison by climate type

In all basins regardless of Climate Type, a general underestimation of the model is seen for the highest runoff months, as noted above looking at the distributions by basin. This is especially evident in Climate Types I and II with pronounced wet seasons as also shown by their lower r-squared values (Figure 4) and lower VE values. Climate Type II also has the highest RMSE value of 4.554 mm/day. Climate Types III and IV have comparable r-squared and VE values though skewness towards underprediction during the highest runoff months is still evident, particularly for Climate Type IV.



Deleted:

Deleted: 2

Deleted: .

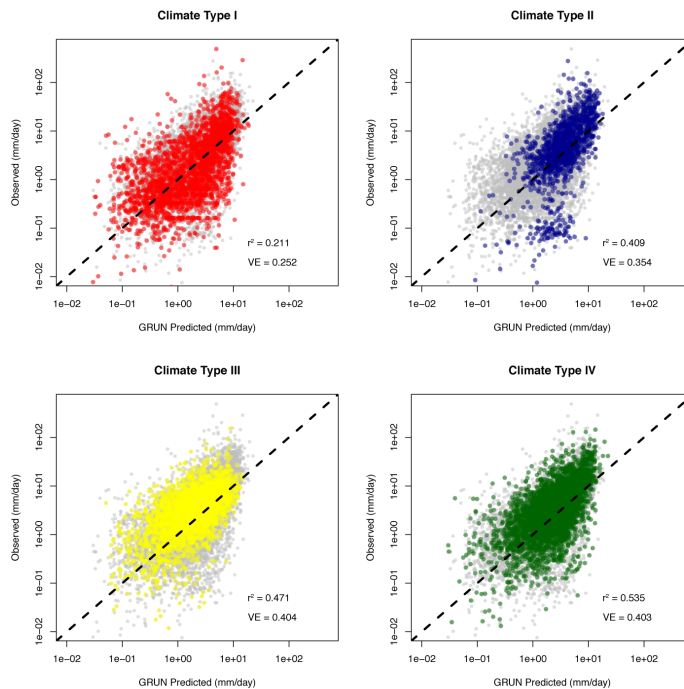
Deleted: colored

Deleted: were

Formatted: Font color: Auto

Deleted: ¶

Deleted: 3



**Figure 4:** Cross plots of GRUN predicted vs. observed monthly runoff by climate type (see Figure 1 for climate type distributions). Grey dots represent all data, colored dots represent data points from that region. The squared pearson correlation coefficient ( $r^2$ ) and volumetric efficiency (VE) metrics are listed on each panel.

Deleted: 3

### 3.3 Correlation and trends with watershed characteristics

In this section we analyse two potential correlations between watershed characteristics and VE scores. In Figure 5 we show a weak positive correlation ( $r^2 = 0.041$ ,  $p = 0.137$ ) between VE and  $\log(\text{catchment size})$  (listed in Table 1) and a stronger negative correlation with mean runoff ( $r^2 = 0.182$ ,  $p < 0.01$ ). However, at low runoff there is significant spread in VE score, driven primarily by Climate Type I catchments (red box and whisker plot in Figure 5C). These catchments experience distinct wet and dry seasons in the northwest Philippines. The positive correlation with catchment size is likely primarily due to the extreme wet months biasing results as described above. This is particularly evident from the Nash-

Deleted: 4

Deleted: 4C



305 Sutcliff Efficiency (NSE) and log(Nash-Sutcliff Efficiency) (NSE-log10) scores in Table 2 and Table A1. Since NSE puts more weight on large flow (Criss and Winston, 2008), it is not surprising that our NSE-log10 scores are in most cases significantly more skilful among our catchments because extremely wet months are weighted less than using the raw runoff values, compared to raw NSE scores. It is also notable that the VE scores using log10 values across the entire dataset is significantly improved (0.363 vs. 0.733; Table 2). The physical significance of these observations are that for large basins the time of concentration of any given flood event will be much longer, thus flood peaks will be wider and subdued due to abstractions and infiltration into the shallow aquifers. This phenomenon is likely less apparent in smaller basins where peak flows are expected to be higher because of less infiltration.

310 Previous studies have investigated the correlation between runoff and catchment size (Mayor et al., 2011), and the different hydrologic and geologic factors that cause non-linear relationships between these two variables (Rodriguez et al., 2014). Recently, Zhang et al. (2019) point out that runoff coefficients increase logarithmically as catchment size decreases. Moreover, the same paper reports that smaller catchments are more sensitive to vegetation cover, slope, and land use compared to larger catchments. This implies that predictability of basin runoff for smaller catchments are increasingly more

315 difficult due to variances in the compounding factors mentioned above. We hypothesize that these effects proposed by Zhang et al. (2019) are also influencing the Philippines streamflow dataset utilized in this study. As such, we suggest that GRUN, is a useful new tool for studying trends, seasonality and average runoff from tropical catchments (such as in previous work: e.g., Merz et al., 2011; Wanders and Wada, 2015) in the Philippines. However, we qualify this finding by noting that GRUN is not suitable for extreme value analyses associated with major tropical storms during the wet seasons unless suitable bias

320 corrections (see next section) can be effectively carried out.

Deleted: As discussed by

Deleted: (

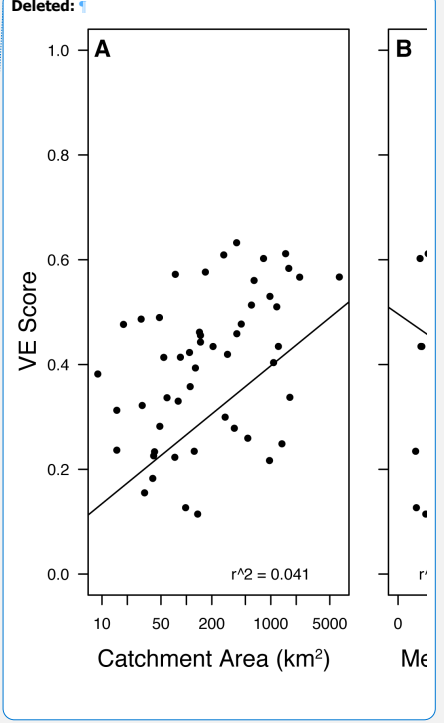
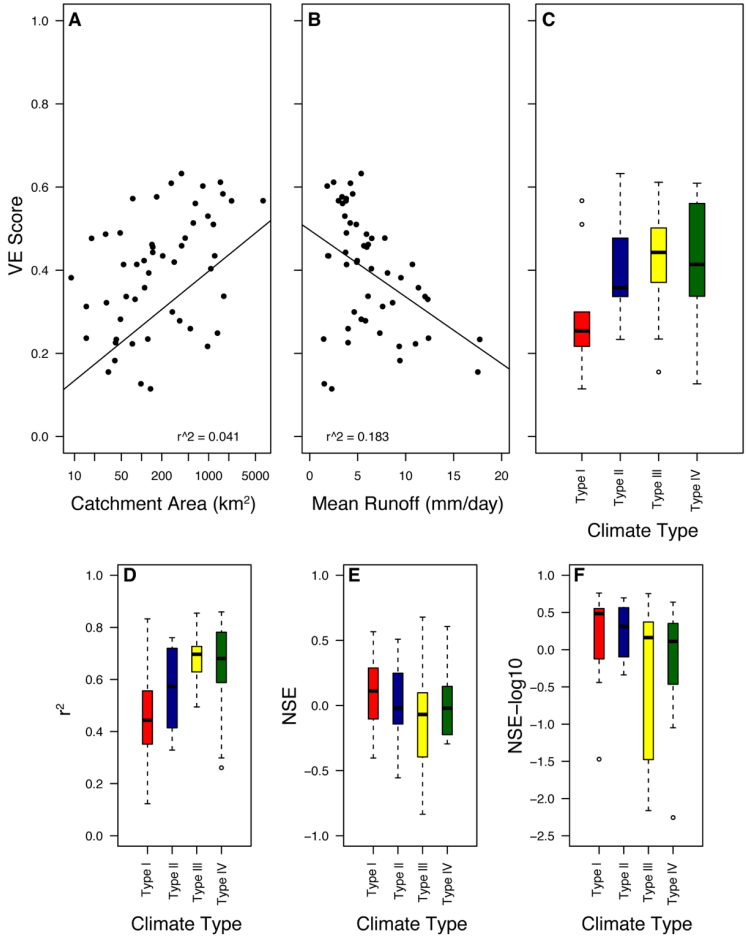
Deleted: ) the

Deleted: .

Deleted: but more importantly that runoff-catchment size relationship derived from large basins must be corrected when applied to smaller catchments..

Deleted: while

Deleted: ,



**Figure 5: Diagnostic plots of volumetric efficiency (VE) results.** Cross plots show the correlation of VE with (A) catchment area and (B) mean runoff. (C) Box and whisker plots show data from distribution of VE by climate type. Box and whisker plots show the median, interquartile range and 95% confidence intervals and outliers (dots). The regression in (A) is between the VE scores and  $\ln(\text{Catchment Area})$ . (D-F) As in (C) for  $r^2$ , NSE and NSE-log10.

### 3.4 Bias Correction and Outlook

Overall, the GRUN data underestimates the actual observed runoff from Philippine basins. The GRUN dataset shows a range of 0 to 10mm/day for most basins and up to 20mm/day for larger basins in the group. The observed maximum runoff values are on average higher by and exceed 50mm/day during extreme rain events (Figure 4). Furthermore, the GRUN dataset also appears to underestimate minimum flow in streams from highly seasonal catchments (e.g. Types I and II).

The underestimation of runoff values during extreme rain events may be a result of the fast saturation of the overlying soil and exceedance of rates of infiltration partly as a result of shallow aquifers filling up and consequently the conversion of excess rainfall into direct runoff (e.g., Tarasova et al., 2018). On the other hand, the underestimation of flow during low flow events may be a result of not accurately accounting for stream baseflow which is fed by shallow aquifers, as well as other effects such as land use and surface properties. These effects may be buffered out in larger catchments leading to our observation described previously that model-data agreement improves with catchment size (Figure 5B).

Given the biases observed in our analysis and in particular the clear underprediction of GRUN during the wettest months we perform a bias correction of the GRUN dataset at a nationwide level using all the available filtered data used in our analysis. We do so in a two-step process to both correct the mean offset and stretch the wettest months to higher values with all transformations occurring in log-transform space (i.e., as displayed in cross plots in Figures 2 and 4). Thus, we first add the mean  $\log_{10}(\text{runoff})$  difference between the observations and the predicted values ( $0.117 \pm 0.022$ ). Following this, using the *lm* function in R, we fit a linear regression between the observations and the GRUN predicted values ( $\log_{10}(\text{runoff, observed}) = m \times \log_{10}(\text{runoff, predicted}) + b$ ) and correct the predicted values using the slope ( $m=0.774 \pm 0.025$ ) and intercept ( $b=0.099 \pm 0.006$ ) derived from this regression. Uncertainties reported here are 68% confidence intervals and were assessed by bootstrap resampling 1,000 observation and prediction pairs without replacement 10,000 times. By carrying out these calculations in log-transform space the highest GRUN runoff values are the most significantly affected, which are the data points we have observed to be most underpredicted in our above analysis (Figures 3A and 4). Because these corrections are being carried out in  $\log_{10}$  space statistical bias (underestimation) is possible (Ferguson, 1986). Following Ferguson (1986) we calculate the unbiased estimate of the variance (s) as 0.0686 mm/day which gives a correction factor (calculated as  $\exp(2.65s^2)$  of 1.0126 may be applied uniformly as a correction factor (multiplier) to the bias corrected values to adjust for possible bias due to the  $\log_{10}$  space regression we have implemented.

To assess this bias correction, we calculated RMSE values at a catchment, climate type and countrywide level (Figure 6 and Tables 2 and A1). The log-transform bias correction greatly improves the nationwide RMSE value by an order of magnitude (2.648 vs. 0.292) and most significantly improves catchments from Climate Types III and IV (Figure 6; 2.285 vs. 0.432 and 2.398 vs. 0.131, respectively; Table 2). Interestingly, the median RMSE value for Climate Type I and II catchments are not significantly improved, however, the RMSE range for both have been reduced (red and blue boxes in Figure 6, respectively).

Deleted: 3

Deleted: 4

These discrepancies may be a function of the direct interpolation of large basin characteristics to smaller basins as discussed in the previous section. The underestimation of runoff values during extreme

Deleted: also

Deleted: .

Deleted: Figure 3).

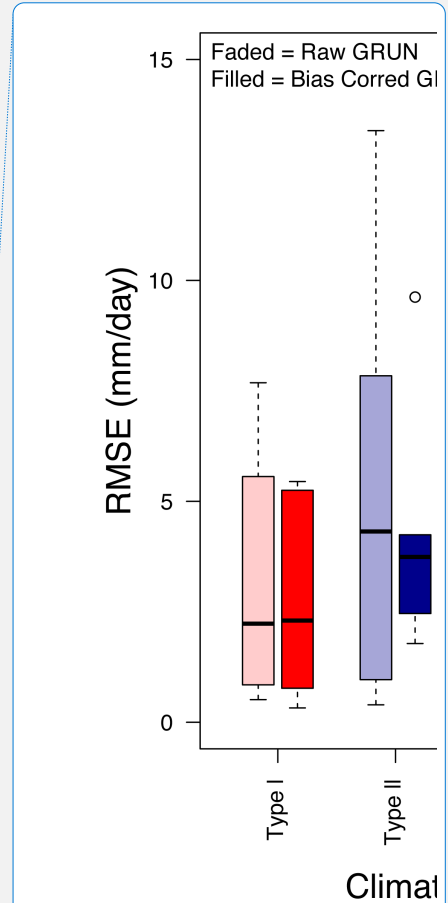
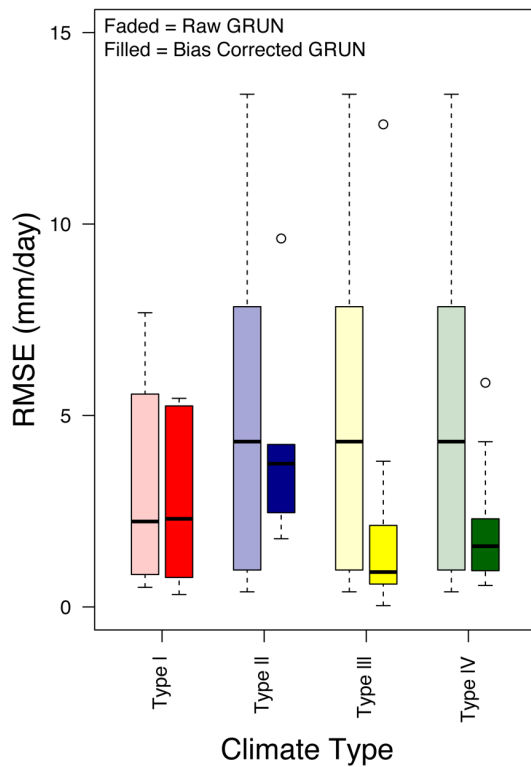
Deleted: 2A and 3).

Deleted: 5

Deleted: 5

Deleted: 5

385 This analysis and the improvement of RMSE values, as well as some of performance metrics such as NSE (see scores tabulated in Table 2), using a simple log-transform based bias correction demonstrates the importance of either: 1) including smaller catchments in future iterations of products such as GRUN, or 2) performing similar bias corrections on a country, region or even catchment scale as appropriate. This is particularly important given that taken at face value the proportional contribution of relatively small tropical land areas to global discharge accounting (e.g., Dai and Trenberth, 2002) would be underestimated without such corrections.



Deleted:  
Deleted: 5  
Deleted: colored  
Deleted: colored

390 **Figure 6.** Root Mean Square Error (RMSE) box and whisker plots of catchments grouped by climate type of observed values versus raw GRUN values (light-colored boxes) and bias-corrected GRUN values (bold-colored boxes). For bias correction equation and country-wide results see Table 2.

400 **4 Conclusion**

Using monthly runoff observations from catchments in the Philippines with more than 10 years of data between 1946 and 2014, we demonstrated across all observations significant but weak correlation ( $r^2 = 0.372$ ) and skillful prediction (Volumetric Efficiency = 0.363 and log(Nash-Sutcliffe Efficiency) = 0.453) between the GRUN-predicted values and actual observations. Looking at different hydrometeorological regimes, we demonstrated that GRUN performs best among low rainfall catchments located in climate types III and IV and showed a weak negative positive correlation between volumetric efficiency and catchment area. Further, we found that particularly for smaller catchments, maximum wet season values are grossly underpredicted by GRUN. The application of a nationwide bias correction to stretch high runoff values using log-transform runoff values greatly improved the RMSE of the predicted values. [Global databases such as GRUN are applicable for aggregated stream discharge estimates and to investigate general trends in the hydrologic characteristics of a region. The recommended bias correction presented here will likely improve such estimates and analysis for the Philippines. While GRUN was never intended to be used for estimating single catchment discharge its applicability for such purposes can be extended provided that proper statistical comparison of modelled versus actual gauged data are initially performed.](#) We thus propose that the utilization of the GRUN dataset can be extended to other ungauged tropical regions with smaller catchments upon applying a similar correction as described in this study.

415

Deleted:

**Acknowledgements**

We thank the Republic of the Philippines Department of Science and Technology Philippine Council for Industry, Energy, and Emerging Technology Research and Development (DOST-PCIEERD) Balik Scientist program for facilitating collaboration between the authors. This project was partially supported through the DOST-PCIEERD project entitled, Catchment Susceptibility to Hydrometeorological Events supporting Tolentino and David, and a DOST-PCIEERD Balik Scientist award to Ibarra. Ibarra is funded by the UC Berkeley Miller Institute and UC President’s Postdoctoral Fellowship. We acknowledge the help of Mart Geronia, Allen Marvin Yutuc, Iris Joy Gonzales, Jia Domingo and EJ Terciano Jr. for help with dataset compilation. DPWH-BRS, GISM and GRDC for datasets availability. [We thank Jaivime Evaristo and an anonymous referee for thorough reviews, and Gionata Ghiggi for ensuring accurate representation of the GRUN product and comments on the manuscript.](#)

**Data availability**

Data was compiled from the DPWH-BRS, GISM and GRDC datasets (see links in text) and is made available as a supplemental file.

**Author contributions**

435            DEI and CPCD designed the study, DEI and PLMT carried out the dataset compilation and screening, PLMT  
verified catchment areas, DEI and PLMT carried out the analysis, and DEI and CPCD prepared the manuscript with  
contributions from PLMT.

**Competing interests**

440    The authors declare that they have no conflict of interest.

## References

- Abon, C. C., David, C. P. C., and Bellejera, N. E. B.: Reconstructing the Tropical Storm Ketsana flood event in Marikina River, Philippines, *Hydrol. Earth Syst. Sci.*, 15, 1283–1289, doi:10.5194/hess-15-1283-2011, 2011.
- 445 Abon, C. C., Kneis, D., Crisologo, I., Bronstert, A., David, C. P. C., Heistermann, M.: Evaluating the potential of radar-based rainfall estimates for streamflow and flood simulations in the Philippines, *Geomat. Nat. Haz. Risk.*, 7, 1390-1405, doi:10.1080/19475705.2015.1058862, 2016
- 450 [Alfieri, L., Lorini, V., Hirpa, F. A., Harrigan, S., Zsoter, E., Prudhomme, C., and Salamon, P.: A global streamflow reanalysis for 1980–2018. \*Journal of Hydrology X\*, 6, 100049, doi: 10.1016/j.hydroa.2019.100049, 2020](#)
- Criss, R. E. and Winston, W. E.: Do Nash values have value? Discussion and alternate proposals. *Hydrol. Process.*, 22(14), 2723-2725, doi:10.1002/hyp.7072, 2008.
- 455 Dai, A., and Trenberth, K. E.: Estimates of Freshwater Discharge from Continents: Latitudinal and Seasonal Variations, *Journal of Hydrometeorology*, 3 (6), 660-687, 2002.
- David, C. P., Cruz, R. V. O., Pulhin, J. M., and Uy, N. M., Freshwater Resources and Their Management, In, Philippine Climate Change Assessment Report WG2: Impacts, Vulnerabilities and Adaptation, OML Foundation, 34-54, 2017.
- 460 [Davie, J.C., Falloon, P.D., Kahana, R., Dankers, R., Betts, R., Portmann, F.T., Wisser, D., Clark, D.B., Ito, A., Masaki, Y., and Nishina, K.: Comparing projections of future changes in runoff from hydrological and biome models in ISI-MIP. \*Earth System Dynamics\*, 4\(2\), 359-374, doi: 10.5194/esd-4-359-2013, 2013](#)
- 465 Do, H. X., Gudmundsson, L., Leonard, M., and Westra, S.: The Global Streamflow Indices and Metadata Archive (GSIM) – Part 1: The production of a daily streamflow archive and metadata, *Earth Syst. Sci. Data*, 10, 765–785, <https://doi.org/10.5194/essd10-765-2018>, 2018.
- Evaristo, J., and McDonnell, J. J.: Global analysis of streamflow response to forest management: *Nature*, 570, 455-461, doi:10.1038/s41586-019-1306-0, 2019.
- 470 [Ferguson, R. I.: River Loads Underestimated by Rating Curves, \*Water Resources Research\*, 22, 74-76, doi: 10.1029/WR022i001p00074, 1986.](#)
- 475 [Ghiggi, G., Humphrey, V., Seneviratne, S.I., Gudmundsson, L. 2019. GRUN: An observations-based global gridded runoff dataset from 1902 to 2014. \*Earth Sys. Sci.\*, 11, 1655-1674, doi:10.5194/essd-11-1655-2019, 2019. Dataset analysed: \[https://figshare.com/articles/GRUN\\\_Global\\\_Runoff\\\_Reconstruction/9228176\]\(https://figshare.com/articles/GRUN\_Global\_Runoff\_Reconstruction/9228176\)](#)
- Gudmundsson, L., Do, H. X., Leonard, M., and Westra, S.: The Global Streamflow Indices and Metadata Archive (GSIM) – Part 2: Quality control, time-series indices and homogeneity assessment, *Earth Syst. Sci. Data*, 10, 787–804, <https://doi.org/10.5194/essd-10-787-2018>, 2018b.
- 480 [Hagemann, S., Chen C., Haerter, J. O., Heinke, J., Gerten, D., and Piani, C.: Impact of a Statistical Bias Correction on the Projected Hydrological Changes Obtained from Three GCMs and Two Hydrology Models, \*Journal of Hydrometeorology\*, 12 \(4\), 556-578, doi: 10.1175/2011JHM1336.1, 2011.](#)
- 485 [Harrigan, S., Zsoter, E., Alfieri, L., Prudhomme, C., Salamon, P., Wetterhall, F., Barnard, C., Cloke, H., and Pappenberger, F.: GloFAS-ERA5 operational global river discharge reanalysis 1979-present. \*Earth System Science Data\*, doi: 10.5194/essd-2019-232, 2020](#)

Deleted:



- Jose, A. M., and Cruz, N. A.: Climate change impacts and responses in the Philippines: water resources, *Clim Res.*, 12 (2-3), 77–84, 1999.
- 495 Kim, H., Watanabe, S., Chang, E. C., Yoshimura, K., Hirabayashi, J., Famiglietti, J., and Oki, T.: Global Soil Wetness Project Phase 3 Atmospheric Boundary Conditions (Experiment 1) [Data set], Data Integration and Analysis System (DIAS), doi:10.20783/DIAS.501, 2017.
- Kintanar R. L.: Climate of the Philippines, PASAGA report, 38, 1984.
- 500 Kumar, P., Masago, Y., Mishra, B. K., and fukushi, K., Evaluating future stress due to combined effects of climate change and rapid urbanization for Pasig Marikina river, Manila, *Groundwater Sustain. Develop.*, 6, 227-34, doi:10.1016/j.gsd.2018.01.004, 2018
- 505 Kumm, M., Guillaume, J. H. A., de Moel, H., Eisner, S., Flörke, M., Porkka, M., Siebert, S., Veldkamp, T. I. E., and Ward, P. J.: The world's road to water scarcity: shortage and stress in the 20th century and pathways towards sustainability, *Nat. Sci. Rep.*, 6, 38495, <https://doi.org/10.1038/srep38495>, 2016.
- 510 Meybeck, M., Kumm, M., and Dürr, H. H.: Global hydrobelts and hydroregions: improved reporting scale for water-related issues, *Hydrol. Earth Syst. Sci.*, 17, 1093–1111, doi: 10.5194/hess-17-1093-2013, 2013.
- Mayor, A. G., Bautista, S., and Bellot, J.: Scale-dependent variation in runoff and sediment yield in a semiarid Mediterranean catchment, *J. Hydrol.*, 397 (1-2), 128-135, doi:10.1016/j.jydrol.2010.11.039.
- 515 Merz, R., Parajka, J., and Blöschl, G.: Time stability of catchment model parameters: Implications for climate impact analyses, *Water Resources Research*, 47 (2), doi: 10.1029/2010WR009505, 2011.
- Nash, J. E. and Sutcliffe, J. V.: River flow forecasting through conceptual models part I – A discussion of principles, *J. Hydrol.*, 10, 282–290, doi:10.1016/0022-1694(70)90255-6, 1970.
- 520 Nurse, L. A., McLean, R. F., Agard, J., Briguglio, L.P., Duvat-Magnan, V., Pelesikoti, N., Tompkins, E. and Webb, A.: Small islands. In: *Climate Change 2014: Impacts, Adaptation, and Vulnerability. Part B: Regional Aspects. Contribution of Working Group II to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change* [Barros, V.R., et al. (eds.)]. Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA, Ch. 29, 1613-1654, 2014.
- 525 Paronda, G. R. A., David, C. P. C., and Apodaca, D. C.: River flow patterns and heavy metals concentrations in Pasig River, Philippines as affected by varying seasons and astronomical tides, *IOP Conf. Series: Earth and Environmental Science* 344, 012049, doi: 10.1088/1755-1315/344/1/012049, 2019.
- 530 Rodríguez, C., Emilio, Cantón, Yolanda, Lazaro, R., et al. 2014. Cross-scale interactions between surface components and rainfall properties. Non-linearities in the hydrological and erosive behavior of semiarid catchments. *J. Hydrol.*. 517:815–825.
- Tarasova, L., Basso, S., Zink, M. and Merz, R.: Exploring Controls on Rainfall-Runoff Events: 1. Time Series-Based Event Separation and Temporal Dynamics of Event Runoff Resposne in Germany, *Water Resources Research*, 54 (10), 7711-7732, doi:10.1029/2018WR022587, 2018.
- 535 Tolentino, P. L. M., poortinga, A., Kanamaru, H., Keesstra, S., Maroulis, J., David, C. P. C., Ritsema, C. J.: Projected Impact of Climate Change on Hydrological Regimes in the Philippines, *PLoS ONE*, 11(10), e0163941, doi: 10.1371/journal.pone.0163941, 2016.
- 540

Wanders, N. and Wada, Y.: Decadal predictability of river discharge with climate oscillations over the 20th and early 21st century, *Geophys. Res. Lett.*, 42, 10689–10695, <https://doi.org/10.1002/2015GL066929>, 2015.

Winsemius, H.C., Aerts, J.C., Van Beek, L.P., Bierkens, M.F., Bouwman, A., Jongman, B., Kwadijk, J.C., Ligtvoet, W., Lucas, P.L., Van Vuuren, D.P. and Ward, P.J.: Global drivers of future river flood risk. *Nature Climate Change*, 6(4), 381-385, doi: 10.1038/nclimate2893, 2016.

WEF: The Global Risks Report 2018, available at: <http://reports.weforum.org/global-risks-2018/>.

Zhang, Q., Liu, J., Yu, X. and chen, L.: Scale effects on runoff and a decomposition analysis of the main driving factors in the Haihe Basin mountainous area. *Sci. Tot. Env.*, 690, 1089-1099, doi: 10.1016/j.scitotenv.2019.06.540, 2019

Table 1: List of stations used in this analysis including full station names, updated catchment areas, years of coverage, division and climate type.

River Name	Station Name	Latitude	Longitude	Coverage	Years of	Catchment	Dataset	Climate Type
					Coverage	Area (km <sup>2</sup> )		
Sinalang River	Penarrubia, Abra	17.622	120.715	1984-2015	32	136.128	BRS	1
Antiqueira River	Sto. Rosario, Antiqueira, Bohol	9.493	123.890	1984-2016	33	54	BRS	4
Amparo River	Brgy. Mabini, Macarohan, So. Leyte	10.042	124.018	1985-2007	23	74	BRS	4
Hirasan River	Upper Hirasan, Rarigara, Leyte	11.258	124.672	1986-2010	27	8.93	BRS	4
Leyte River	San Joaquin, Capocan, Leyte	11.880	124.829	1985-2007	23	29.15	BRS	3
Surigao River	Surigao City	9.796	125.808	1986-2010	25	85	BRS	4
Bais River	Cabunlatan, Bais City, Negros Oriental	9.876	124.140	1989-2015	27	41	BRS	2
Lingayon River	Lingayon, Alang-Alang, Leyte	11.192	124.863	1957-1991	35	18	BRS	4
Sapiniton River	Libton, San Miguel, Leyte	11.188	124.795	1984-2010	27	277.3	BRS	4
Laosag River	Poblacion, Laosag City, Ilocos Norte	18.203	120.590	1984-2016	33	1355	BRS; updated from Tolentino et al. (2016)	1
Pared River	Baybayog, Alcala, Cagayan	17.682	121.270	1983-1996	14	966	BRS; used in Tolentino et al. (2016)	1
Gumano River	Ipil, Echague, Isabela	16.812	121.211	1986-2001	16	977	BRS; used in Tolentino et al. (2016)	3
Magat River	Baretel, Baguio, Nueva Vizcaya	16.992	121.073	1986-2002	17	2199	BRS; used in Tolentino et al. (2016)	3
Camiling River	Poblacion, Mayantoc, Tarlac	15.018	120.503	1985-2017	33	288	BRS; updated from Tolentino et al. (2016)	1
Gumain River	Sta. Cruz, Lubao, Pampanga	14.960	120.441	1985-2001	17	370	BRS; used in Tolentino et al. (2016)	1
Rio Chico River	Sto. Rosario, Zaragoza, Nueva Ecija	15.658	120.088	1985-2006	22	1177	BRS; used in Tolentino et al. (2016)	1
San Juan River	Porac, Calamba, Laguna	14.498	121.267	1986-1999	14	165	BRS; used in Tolentino et al. (2016)	4
Pangalan River	Pangalan, Pinamlayan, Oriental Mindoro	13.275	121.260	1989-1999	11	32	BRS; used in Tolentino et al. (2016)	3
Das-ay River	Sto. Nino II, Hinunangan, Leyte	10.385	125.202	1987-2007	21	59	BRS; used in Tolentino et al. (2016)	2
Tukuran River	Tinotongan, Tukuran, Zamboanga del Sur	7.627	123.030	1986-2009	24	147	BRS; used in Tolentino et al. (2016)	3
Hijo River	Apokan, Tagum, Davao del Norte	7.812	125.211	1986-2016	31	634	BRS; updated from Tolentino et al. (2016)	4
Cagayan de Oro River	Cabula, Cagayan de Oro City, Misamis Oriental	8.316	124.811	1991-2004	14	1079	BRS; used in Tolentino et al. (2016)	4
Davao River	Tagatito, Davao City	7.329	125.634	1984-1999	16	1683	BRS; used in Tolentino et al. (2016)	4
Alili River	Impos, Iloilo, Sultan Kudarat	6.568	124.085	1980-1994	15	1231	BRS; used in Tolentino et al. (2016)	3
Agusan Canyon River	Camp Phillips, Manolo Fortich, Bukidnon	8.296	124.450	1986-2004	19	48	BRS; updated from Tolentino et al. (2016)	3
Wawa River	Wawa, Bayugan, Agusan del Sur	8.261	125.501	1981-2010	30	396	BRS; updated from Tolentino et al. (2016)	4
Buyan River	Malandag, Malungon, South Cotabato	6.317	125.749	1986-2004	19	207	BRS; used in Tolentino et al. (2016)	4
Gaugas River	Manalpac, Solosona	18.080	120.830	1978-1988	11	73	GISM	1
Jalaur River	Calyan, Pototan, Iloilo	10.930	122.670	1976-1988	13	1499	GISM and GRDC	3
Padian River	Bangay	18.080	120.700	1946-1979	34	534	GRDC	1
Pampanga River	San Agustin	15.170	120.780	1946-1977	32	6487	GRDC	1
Sipocot River	Sabang	13.810	122.990	1946-1970	25	447	GRDC	2
Mambasao River	Tumalalad	11.260	122.570	1950-1978	29	307	GRDC	3
Padula River	Lapulao	6.660	125.280	1949-1978	30	821	GRDC	4
Aloran River	Juan Baccay, Aloran, Misamis Occ.	8.420	123.820	1978-2003	26	30	GRDC + BRS	3
Cabacan River	Baduang, Pagadpad	18.580	120.800	1979-2017	39	60	GRDC + BRS	1
Maragayap River	Sta. Rita, Bacnotan, La Union	16.750	120.374	2004-2017	14	40	BRS	1
Abacan River	San Juan, Mexico, Pampanga	15.118	120.703	2004-2017	14	217	BRS	1
Hibayog River	La Victoria, Carmen, Bohol	9.876	124.141	2004-2017	14	41	BRS	4
Manaba River	Calma, Garcia-Hernandez, Bohol	9.631	124.131	2001-2016	16	98	BRS	4
Gubayan River	Canawa, Candijay, Bohol	9.848	124.450	2001-2017	17	48.5	BRS	4
Bangkrohan River	Brgy. Tagaytay, Bato, Leyte	10.342	124.834	1984-1996; 2000-2009	17	168	BRS	4
Borongan River	Brgy. San Mateo, Borongan City	11.628	125.403	1990-2008	19	111	BRS	2
Loom River	Brgy. Calico-an, Borongan City	10.594	125.404	1986-2004	19	42	BRS	2
Pagbanganan River	Brgy. Makinhas, Baybay City	10.637	124.865	1984-2008	25	128	BRS	4
Rizal River	Brgy. Rizal, Babatngon, Leyte	11.389	124.908	1990-2008	18	15	BRS	4
Tenani River	Brgy. Tenani, Paranas (Wright), Samar	11.806	125.127	1985-2001	17	394	BRS	2
Diakan River	Diakan, Mamukan, Zamboanga del Norte	8.480	123.048	1985-1991; 1997-2000	11	109	BRS	3
Kabasalan River	Banker, Kabasalan, Sibugay, Province	7.831	122.778	2002-2011	10	143	BRS	3
Sindangan River	Diocoyong, Sindangan, Zamboanga del Norte	8.217	123.057	1990-2003	14	590.5	BRS	3
Alahijid River	Alahijid, Misamis Oriental	8.570	124.476	1991-2009	19	124	BRS	3
Kipalaja River	Tiburcia, Kapalong, Davao del Norte	7.602	125.681	2004-2016 1987-1995;	13	147	BRS	4
Banaue River	Poblacion, Banaue, Ifugao	16.915	121.061	2005-2010	15	15	BRS	3
Aciga River	Santiago, Agusan del Norte	9.269	125.570	2002-2015 1982; 1984- 1987; 1989-	14	80	BRS	4
Agusan River	Sta. Josefa, Agusan del Sur	7.993	126.036	2010	27	1633	BRS	4

Table 2: Results of statistical agreement between GRUN aggregated by Climate Type and for the entire dataset (see Table A1 for individual catchments)

	Pearsons Coeff ( $r^2$ ) *	Volumetric Efficiency (VE)	Nash-Sutcliffe Efficiency (NSE)	Nash-Sutcliffe Efficiency (NSE- log10)	Root Mean Square Error	Root Mean Square Error Bias Corrected	Volumetric Efficiency (VE) Bias Corrected **	Nash-Sutcliffe Efficiency (NSE) Bias Corrected **	Nash-Sutcliffe Efficiency (NSE- log10) Bias Corrected **
Entire Dataset	0.372	0.363	0.091	0.453	2.648	0.292	0.323	0.182	0.385
Entire Dataset log10(runoff)	0.546	0.733	n/a	n/a	n/a	n/a	1.067	n/a	n/a
Climate Type 1 (n=12)	0.211	0.252	0.062	0.538	2.476	0.298	0.168	0.111	0.432
Climate Type 2 (n=6)	0.409	0.354	0.05	0.49	4.554	0.544	0.349	0.188	0.457
Climate Type 3 (n=15)	0.471	0.404	0.026	0.23	2.285	0.432	0.345	0.011	0.188
Climate Type 4 (n=22)	0.535	0.403	0.159	0.414	2.398	0.131	0.377	0.323	0.36

Notes

\* For regressions forced through intercept of 0

\*\* Two-step bias correction procedure where first mean offset is added to the predicted GRUN values and then a log-transform stretch correction is applied (see text for details)