

Author's Response hess-2014-240: Evaluation of the satellite-based Global Flood Detection System for measuring river discharge: Influence of local factors

B. Revilla-Romero, J. Thielen, P. Salamon, T. De Groeve, and G.R. Brakenridge

Author Comments #1:

We would like to thank Reviewer#1 for their careful consideration of this manuscript (hess-2014-240) and for their helpful and insightful comments. We have carefully considered the reviewer's comments and worked to include them in the revised version of the manuscript according to the proposed suggestions.

Please find below the responses to the reviewer's comments.

General comments

This paper is an interesting evaluation of the skill of the Global Flood Detection System to measure river discharge from satellite passive microwave signals, and is certainly worthy of publication after some correction. The correlation between the daily ground station-measured water discharge and the satellite signal is measured for a range of rivers with different widths, floodplain areas, land cover types, climatic regions and other factors. For African, Asian and North American rivers, the mean R values are less than 0.5, and the correlation is only medium. Only European and South American rivers give high correlation (>0.5). It might be argued that a judicious set of ranges of R has been employed ($R < 0.3$, $0.3 - 0.7$, >0.7), in which many rivers lie in the middle range, but still may have R values < 0.5 , so that the correlation is only medium. The authors should comment on this. The relatively low R values show the difficulty of obtaining a reasonable signal-to-noise ratio from a 10km pixel when the flood width is often substantially less than the pixel width. As a result, it is obviously a sensible idea to identify sites where the method will work because of the associated site variables, and use these for future studies, rather than trying to make the method work for all sites. The method would also appear to work best for detecting floods rather than forecasting them, since a 4-day average signal is used, partly to cope with the time lag between changes in stage at a gauging station and associated changes in flood extent.

Specific comments

7333/14: Make it clear that you are talking about river floods (or does this include deaths in the tsunami of 2001?).

Author's reply: Modified in the manuscript as suggested by the Reviewer.

7337/6 In a flood situation, is the error on the observed discharge not higher than the 5–20% quoted?

Author's reply: Explanation added on the manuscript as suggested.

The uncertainty of river discharge is higher during floods events when the stage-discharge relationship, the so-called rating curve, is used. As evaluated by Pappenberger et al. (2006), the analysis of rating curve uncertainties leads to an uncertainty of the input of 18–25% at peak discharge. Di Baldassarre and Montanari (2009) showed that the total rating curve errors increase, when the river discharge increases and varies from 1.8% to 38.4% with a mean value of 21.2%.

7342/13 What was the spread of the R2 values for the fits?

Author's reply: answered below.

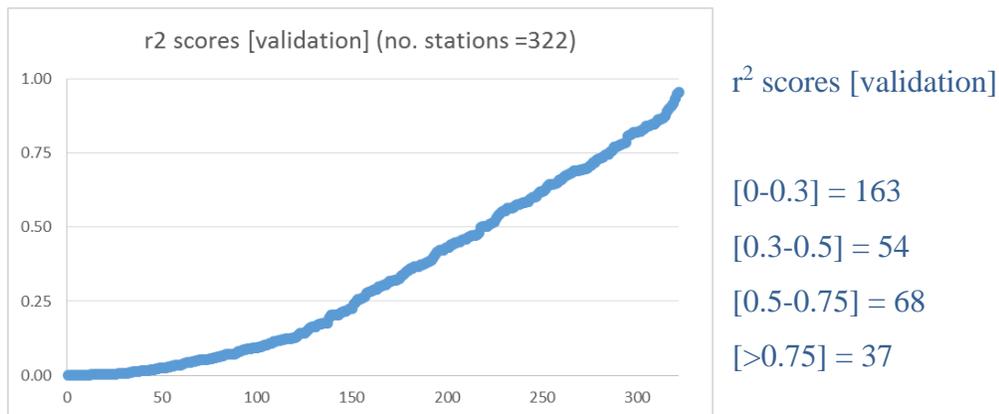


Figure 1. r² scores obtained on the validation.

In Fig. 3b, please make it clearer that different rating equations are being used for different months, not simply that in fig. 3a.

Author's reply: Clarified in the figure caption as suggested by the reviewer.

In fig. 3a, why aren't there 15 points on the graph, one for each March between 1998-2002?

Author's reply: The calibration was done for 5 years (1998-2002), therefore the five points represent the five mean values for March in this case.

7344/20 A little more description of the **Gini index** might help the reader.

Author's reply: explanation added on the manuscript.

Gini's mean difference was first introduced by Corrado Gini in 1912 as an alternative measure of variability and the parameters derived from it, such as the Gini index, also referred to as the concentration ratio (Yitzhaki and Schechtman, 2013). The Gini index is mostly popular in economics, however it is also used in other areas, such as building decision trees in statistics to measure the purity of possible child nodes, and it has been compared with other equality measures (Gonzales,L., et al. 2010).

How does the **Random Forest** method cope if the variables are correlated (as e.g. discharge and river width probably are)? Is the correlation between variables output from the method as they would be from a principal component analysis? If so, it would be useful in the subsequent analysis to know the correlations between variables to know which were most significant.

Author's reply: explanation added on the manuscript.

The random forests algorithm, introduced by Breiman (2001), is a modification of bagging that aggregates a large collection of tree-based estimators and has better estimation performances than a single random tree: each tree estimator has low bias but high variance whereas the aggregation achieve a bias-variance trade-off. This algorithm has good predictive performances in practice, they work well for high dimensional problems and they can be used with multi-class output, categorical predictors and imbalanced problems. Moreover, the random forests provide some measures of the importance of the variables with respect to the prediction of the outcome variable. The random forest algorithm was selected instead of the Principal Component Analysis as we had mixed data types because some of the variables to be study were categorical instead of continuous.

Although, the effect of the correlations on these measures has been studied recently (see Archer and Kimes (2008), Strobl et al. (2008), Nicodemus et al. (2010), Nicodemus (2011), Auret and Aldrich (2011), Tolosi and Lengauer (2011), Grömping, U. (2009) and Gregorutti et al. (2013)) there is no yet a consensus on the interpretation of the importance measures when the predictors are correlated and on what is the effect of this correlation on the importance measure.

In order to test the effect on the results when correlated variables were included in the analysis, an independent Random Forest analysis was carried out during the analysis (not shown in the paper) for the same variables but excluding the river width and the presence of floodplains and wetlands variables. Results also showed that the mean daily observed discharge had the highest importance and the presence of hydraulic structures (mainly dams) and of river ice had the lowest importance to classify a location as good or poor performance.

7345/25 Do you really mean that the signal may have a large natural variation, or that the noise is instrument noise?

Author's reply: answered and edited on the manuscript as follows.

We meant that the signal to noise ratio might be low for a site or have intermittent instrument noise occasionally producing intermittent positive spikes in discharge. We have edit this in the manuscript.

73446/8 $R = 0.3$ is chosen as a threshold in fig. 4, yet this is only a medium correlation. What happens if you chose $R = 0.5$ as the threshold, are there too few sites satisfying this criterion then?

Author's reply: For this study, 42 sites have $R > 0.5$

7346/23 In fig. 5, in the eastern USA, many stations had $R > 0.3$ in the calibration (fig. 4), but have $NSE < 0$ in validation. Why is this? The rivers are presumably often wide and on floodplains near the sea at these observation points?

Author's reply: explanation added on the manuscript. Not the map below as it is complementary to figures already shown in the manuscript.

Figure 5 doesn't shows the calibration score. It shows the initial correlation between GFDS signal vs. in situ observed discharge.

The figure below (Fig. 2 of the Author's comments) shows the R score obtained during the validation for stations located in Eastern USA (no. of stations=66). In addition, it shows the pixels values when the river width is higher than 1km (Yamazaki et al., 2014) and the Global Lake and Wetland Database layers (Lehner and Doll, 2014).

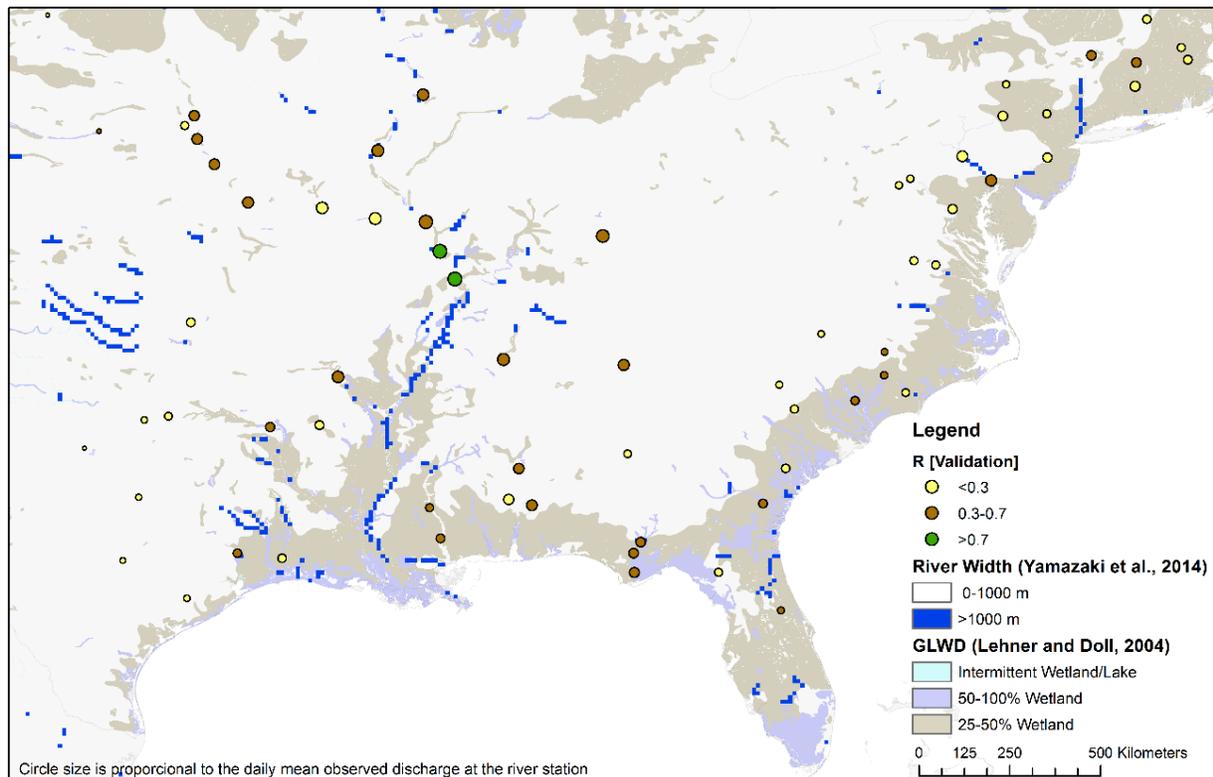


Figure 2. R score of the validation (n=66 station) for Eastern USA.

As shown in the manuscript for the whole of the stations, we conclude that most of the stations in this region obtained poor scores due to a number of factors: ~64% of these stations have a mean discharge value lower than $500 \text{ m}^3\text{s}^{-1}$ and ~88% of the stations are located at river width lower than 1km. In addition, ~59% of the stations are located in wetlands areas. Sites with these characteristics might not provide useful outputs when aiming to measure river discharge through the use of satellite flood signal, as it is the case of some of this stations.

7347/10 It is probably true that locations with a river width higher than 1 km are more likely to score an R larger than 0.3, but it would be worth quantifying R for widths > 1km and showing that it's significantly larger than 0.3. Author's reply: Quantification added on the manuscript. The mean R score is 0.60. Where 26 out of 64 (~41%) have R > 0.75.

A related point is, in fig. 6a, could you explain why some rivers of 100m or less width have R values as high as the widest rivers? Intuitively you would have thought the brightness temperature for a pixel containing water would depend on the river width (perhaps I'm confusing the river width with the flood width here?).

Author's reply: explanation added to the manuscript.

The retrieval of the satellite signal also depends on the floodplain geometry. As soon as the river floods and water goes over-bank, the proportion of water in the wet pixel greatly increases. So the score should be also high for small rivers with a proportionally big floodplain.

7347/24 might not provide reliable results: : It would be better to quantify this rather than just stating it. You could use a statistical test to compare the rivers with $Q < 500 \text{ m}^3/\text{s}$ that have $R < 0.3$ with rivers with $Q > 500$ that have $R > 0.3$, and show that they were significantly different.

Author's reply: Quantification added to the manuscript.

As 77% of the stations with $Q < 500 \text{ m}^3/\text{s}$, have $R < 0.3$, while 91.5% of the stations with $Q > 500 \text{ m}^3/\text{s}$ have $R > 0.3$, locations with discharge of less than $500 \text{ m}^3/\text{s}^{-1}$ might not provide reliable results for a global satellite-based monitoring system.

Technical corrections

All comments were adapted according to the Reviewer's suggestion.

7333/16 Golnaraghi 2009 and Kundzewicz 2012 refs missing /28 UNOSAT 2013 ref missing.

Modified

7335/20 climate-drive -> climate-driven /27 global -> a global 7337/9 us -> as

Modified

7340/4 Example -> Examples /17 define M/C signals /22 split sentence at 'an array'

Modified

7345/8 as validated -> were validated /13 calibrate -> calibrated /14 discharge satellites -> satellite discharge

Modified

7346/16 two-years -> two years /20 shorted -> shorter

Modified

7348/8 25x 25 pixel -> 25 x 25 km pixel /28 To note -> Note

Modified

7349/22 Where highest -> The highest

Modified

7350/8 presence or not -> presence or absence /20 for - the most of – the -> for most of the

Modified

7351/12 in some -> on some

Modified

7352/2 test -> tested

Modified

7354/ fig 12 caption: was chose -> was chosen; of the stations -> or the stations; station -> stations

Modified

7353/2 replace the semicolons with commas in this long sentence /20 satellite measured -> satellite-measured

Modified

7354/10 no verb in sentence /15 a more -> more

Modified

Author Comments #2:

We would like to thank Dr Guy Schumann for his very useful and constructive comments on the paper (hess-2014-240). We have carefully considered the reviewer's comments and worked to include them, when considered appropriate, in the revised version of the manuscript according to the proposed suggestions. Please find below the responses to the reviewer's comments.

Review of the paper 'Evaluation of GFDS' by B. Revilla-Romero et al.. This is an interesting paper reviewing the factors influencing the accuracy of discharge measurement as provided by GFDS. Papers of this type (i.e. evaluation of global Earth monitoring systems and identification/discussion of influencing factors) are highly valuable and absolutely necessary to add both scientific credibility and reliability to a global measurement or/and model system, which will ultimately lead to an increased fidelity and 'trust' in that system by the end-user/decision.

In my opinion, this paper should be published in HESS after addressing some minor to moderate comments:

- Introduction (7334, top of page): Please mention also the International Disaster Charter and efforts such as CEOS etc. in view of space-based support of relief services during disasters.

Author's reply: citation completed on the manuscript as suggested by reviewer.

- 7335 (L 26): Replace 'compared' with 'comparable'.

Author's reply: modified on the manuscript.

- 7337 (L 14): I agree with this statement but would appreciate if the authors added a sentence to this along the lines of: 'the extent to which this is true needs to be fully investigated however.'

Author's reply: added to the manuscript as suggested by reviewer

- 7342: I understand that you want to use linear equations for simplicity but would a simple power-law function not yield similar or better result. Have the authors tried that?

Author's reply:

To test the results using this suggested alternative approach, we used a power-law function:

$$y = k * x^n \quad (1)$$

where y is the in situ observed discharge and x the satellite signal. Then, taking the logarithm of both sides of the Eq.1 yields the linear equation:

$$\log_{10} y = \log_{10} k + n * \log_{10} x \quad (2)$$

where $\log_{10} k$ is the intercept and n the slope. After the calculation of both constants k and n , the power law function can be used to calibrate the GFDS signal into discharge units (m^3s^{-1}) as done when using the linear regression approach. An example is shown in Figure 3c of this document for the same station shown on Figure 3 of the manuscript (Senanga: Long 23.25, Lat. -16.116; Zambezi River).

Applying the different power law functions obtained for each month to the GFDS signal for the same two-year period as on the manuscript, alike GFDS measured discharge values were obtained. The skill scores achieved for the validation using the power law (e.g. Fig. 3d) are similar to those obtained using linear regression. In view of the results and although this approach also produce valid results, we prefer to leave the methodology as it is on the manuscript.

Added to the manuscript: (section 3.2) Power law fitting was also tested to calibrate the signal into discharge units yielding similar results (see Open Discussion Author's Response).

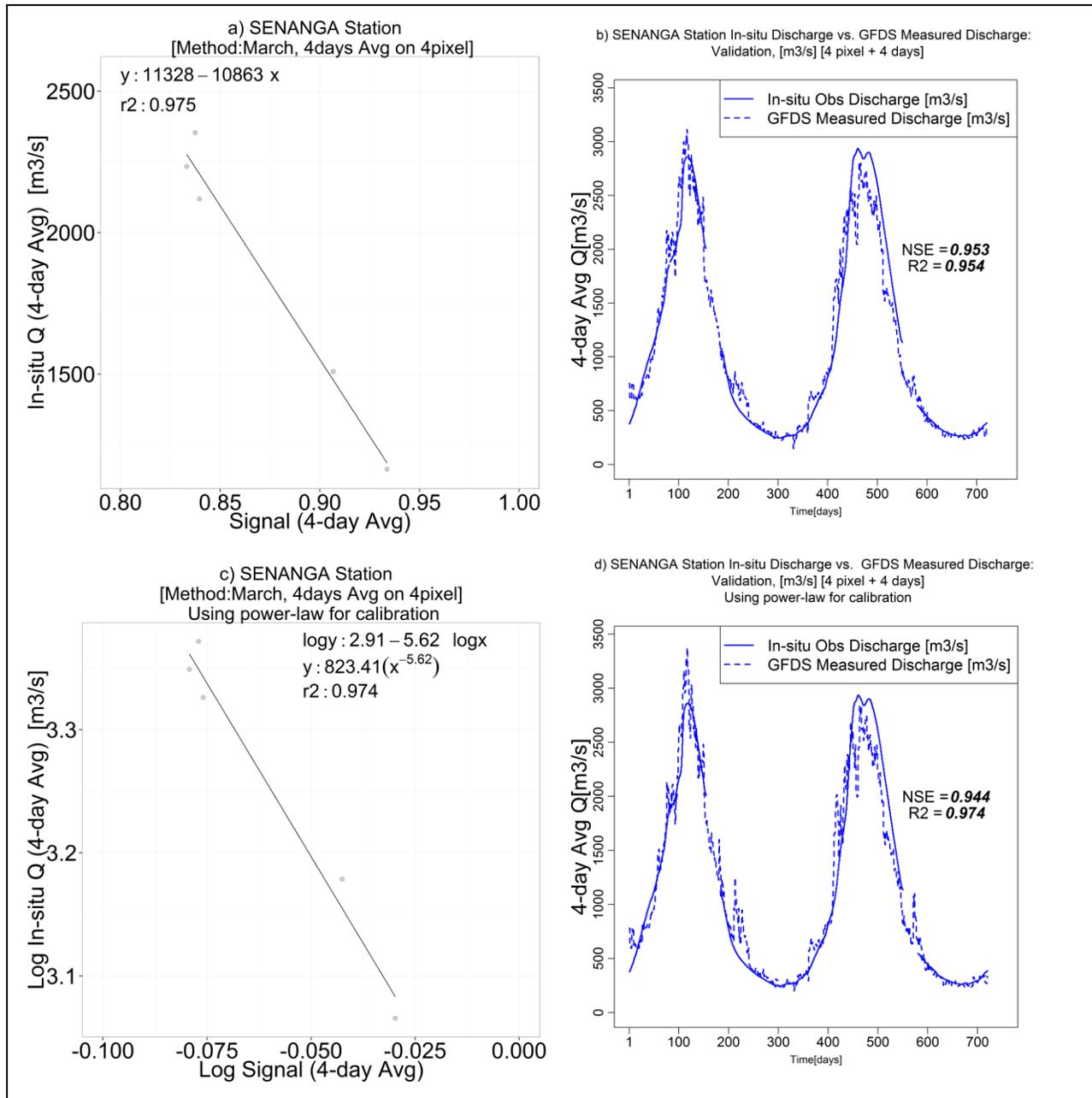


Figure 3 (a) Scatterplot for the Senanga station (Long 23.25, Lat. -16.116) in the Zambezi River (Africa). (b) Validation hydrograph for 2003–2004 and skill scores for Senanga. The (monthly) linear rating equation was used to calibrate the signal into discharge units. Different rating equations were used for different months. (c) Scatterplot for the Senanga station. (d) Validation hydrograph for 2003–2004 and skill scores for Senanga. The (monthly) power-law function was used to calibrate the signal into discharge units. Different equations were used for different months. Note that Figure 3a of the manuscript was amended due to a typo mistake on the linear equation values. (Fig. 3 (a,b) of this document correspond to Fig. 3(a,b) on the manuscript)

- 7343: Since there may be a non-linearity between the station Q and the satellite Q as argued on the previous pages (time lag, etc.), why not employ a Spearman correlation? The Pearson assumes linearity.

Author's reply:

By using a rating equation for each month individually, instead of a single rating equation for the full period, to calibrate the signal into discharge units the derived daily discharge values adjusted better on the timing and also during low flow periods.

As suggested by reviewer, Spearman correlations were calculated for all the stations. Table 1 shows the continental average Pearson and Spearman skill scores. For all continents, average higher values (~106%) were obtained using Spearman correlation in comparison with Pearson score. However, we argue that changing the skill score used on this part of the analysis and for illustration on figures 6-12, will not impact the main findings presented on this manuscript.

Table 1. Mean continental Pearson and Spearman skill scores, obtained on the validation.

*4 stations were excluded on the calculation due to accidental data loss, therefore score varies from manuscript.

Continent	Mean Pearson	Mean Spearman
Africa	0.382	0.403
Asia	0.358	0.438
Europe	0.508	0.537
North America	0.502*	0.538
South America	0.694	0.720
Total	0.527	0.560

Added to the manuscript: (section 3.3) Spearman's rank correlation coefficient was also calculated to assess the validation performance. While Pearson benchmarks linear relationship, Spearman benchmarks monotonic relationship. Spearman's validation scores just obtained a mean value 6% higher than Pearson mean score (see Open Discussion Author's Response). On this manuscript, results are analysed based on the scores obtained using Pearson correlation coefficient.

- 7343: The NSE as argued is showing skill in some data or model when $NSE > 0$ since $NSE = 0$ means as good as mean in observed data, so why not consider the fact that when $NSE > 0$, then the use of satellite discharge should be preferred to long term observed mean, which means 'satisfactory' but not 'good' performance.

Author's reply: Sentence edited on the manuscript. Results and figures were already showed number of stations with $NSE > 0$ and $NSE > 0.50$

- Of course the completeness or incompleteness of each discussion section about the factors influencing the validation / calibration results can be argued forever but I think as a first step analysis and discussing the main factors these sections give a very good appreciation. For that reason maybe the title could be changed to: '.... : a first analysis of the influence of local factors', but I leave that decision to the authors and editor(s).

Author's reply: We acknowledge this suggestion, but would prefer not to make the manuscript's title longer.

- It is great that there is a lot of future work planned on this topic by the team – looking forward to it.

Author's reply: Thank you for your encouragement.

References of the author's reply section.

- Archer, K. J. and Kimes, R. V. Empirical characterization of random forest variable importance measures. *Computational Statistics and Data Analysis*, 52:2249–2260, 2008. doi: 10.1016/j.csda.2007.08.015
- Auret, L. and Aldrich, C. Empirical comparison of tree ensemble variable importance measures. *Chemometrics and Intelligent Laboratory Systems*, 105:157–170, 2011. doi: 10.1016/j.chemolab.2010.12.004
- Breiman, L.: Random forests, *Mach. Learn.*, 45, 5–32, 2001.
- Committee on Earth Observation Satellites (CEOS) Flood Pilot, <http://www.ceos.org/>, last access: 1 September 2014.
- Di Baldassarre, G. and Montanari, A.: Uncertainty in river discharge observations: a quantitative analysis, *Hydrol. Earth Syst. Sci.*, 13, 913–921, doi: 10.5194/hess-13-913-2009, 2009.
- Disaster Charter, 2013. Space and Major Disasters. <http://www.disasterscharter.org/>, last accessed 1 September 2014.
- Gonzalez, L., Velasco Morente, F., Gavilan Ruiz, J.M., Sanchez-Reyes Fernandez, J.M. The Similarity between the Square of the Coefficient of Variation and the Gini Index of a General Random Variable. *Journal of Quantitative Methods for Economics and Business Administration* 10: 5–18.2010, ISSN 1886-516X.
- Gregorutti, B., Michel, B., Saint-Pierre, P. Correlation and variable importance in random forests. Cornell University Library, 2013. arXiv: 1310.5726 [stat]
- Grömping, U. Variable Importance Assessment in Regression: Linear Regression versus Random Forest. *The American Statistician*. 11/2009; 63:308–319, 2009. doi: 10.1198/tast.2009.08199
- Lehner, B. and Döll, P.: Development and validation of a global database of lakes, reservoirs and wetlands, *J. Hydrol.*, 296/1–4, 1–22, 2004. doi: 10.1016/j.jhydrol.2004.03.028
- Nicodemus, K. K. Letter to the editor: On the stability and ranking of predictors from random forest variable importance measures. *Briefings in Bioinformatics*, 12:369–373, 2011. doi: 10.1093/bib/bbr016

- Nicodemus, K. K., Malley, J. D., Strobl, C., and Ziegler, A. The behaviour of random forest permutation-based variable importance measures under predictor correlation. *BMC Bioinformatics*, 11:110, 2010. doi: 10.1186/1471-2105-11-110
- Pappenberger, F., Matgen, P., Beven, K.J., Henry, J., Pfister, L., Fraipont de, P., Influence of uncertain boundary conditions and model structure on flood inundation predictions, *Advances in Water Resources*, Volume 29, Issue 10, Pages 1430-1449, 2006. doi: 10.1016/j.advwatres.2005.11.012
- Strobl, C., Boulesteix, A.-L., Kneib, T., Augustin, T., and Zeileis, A. Conditional variable importance for random forests. *BMC Bioinformatics*, 9:307, 2008. doi: 10.1186/1471-2105-9-307
- Tolosi, L. and Lengauer, T. Classification with correlated features: unreliability of feature ranking and solutions. *Bioinformatics*, 27:1986–1994, 2011. doi: 10.1093/bioinformatics/btr300
- Yamazaki, D., O’Loughlin, F., Trigg, M. A., Miller, Z. F., Pavelsky, T. M., and Bates, P. D.: Development of the global width database for large rivers, *Water Resour. Res.*, 50, 3467–3480, doi: 10.1002/2013WR014664, 2014.
- Yitzhaki, S., Schechtman, E. *The Gini Methodology. A Primer on a Statistical Methodology*. 2013. Springer Series in Statistics. Volume 272, 2013, ISBN: 978-1-4614-4720-7.

LIST OF ALL CHANGES MADE IN THE MANUSCRIPT:

- Text: all modifications where it was affirmed “added/edited/modified on the manuscript” on the present author’s replies section.
- Tables: additional tables were not added.
- Figures: Fig.3a due to a typo mistake on the original figure.
- References: added all new references from the below list.

Please find below a marked-up manuscript version with the modifications made on the manuscript (highlighted in both red and blue)

1 Evaluation of the satellite-based Global Flood Detection 2 System for measuring river discharge: Influence of local 3 factors

4 **B. Revilla-Romero**^{1,2}, **J. Thielen**¹, **P. Salamon**¹, **T. De Groeve**¹, **G. R. Brakenridge**³

5 [1] European Commission Joint Research Centre, Ispra, Italy

6 [2] Utrecht University, Faculty of Geosciences, Utrecht, the Netherlands

7 [3] University of Colorado, Boulder, USA

8 Correspondence to: B. Revilla-Romero (beatriz.revilla-romero@jrc.ec.europa.eu)

9 **Abstract**

10 One of the main challenges for global hydrological modelling is the limited availability of
11 observational data for calibration and model verification. This is particularly the case for real time
12 applications. This problem could potentially be overcome if discharge measurements based on
13 satellite data were sufficiently accurate to substitute for ground-based measurements. The aim of
14 this study is to test the potentials and constraints of the remote sensing signal of the Global Flood
15 Detection System for converting the flood detection signal into river discharge values.

16 The study uses data for 322 river measurement locations in Africa, Asia, Europe, North America
17 and South America. Satellite discharge measurements were calibrated for these sites and a
18 validation analysis with in situ discharge was performed. The locations with very good
19 performance will be used in a future project where satellite discharge measurements are obtained
20 on a daily basis to fill the gaps where real time ground observations are not available. These include
21 several international river locations in Africa: Niger, Volta and Zambezi rivers.

22 Analysis of the potential factors affecting the satellite signal was based on a classification decision
23 tree (Random Forest) and showed that mean discharge, climatic region, land cover and upstream
24 catchment area are the dominant variables which determine good or poor performance of the
25 measurement sites. In general terms, higher skill scores were obtained for locations with one or
26 more of the following characteristics: a river width higher than 1km₂, a large floodplain area and
27 in flooded forest₂ with a potential flooded area greater than 40%₂, sparse vegetation, croplands or
28 grasslands and closed to open and open forest₂, Leaf Area Index > 2₂, tropical climatic area₂ and

29 without hydraulic infrastructures. Also, locations where river ice cover is seasonally present
30 obtained higher skill scores. The work provides guidance on the best locations and limitations for
31 estimating discharge values from these daily satellite signals.

32 **Keywords:** Floods; Passive microwave sensing; Discharge measurement; global evaluation; local
33

34 **1 Introduction**

35

36 Flooding is the most prevalent natural hazard at the global scale, often with dire humanitarian and
37 economic effects. According to the International Disaster Database (EM-DAT), an average of 175
38 flood events per year occurred globally between 2002-2011, affecting an average of 116.5 million
39 people, and causing economic losses of US\$25.5 billion. According to MunichRe (2014), the
40 costliest natural catastrophe worldwide in terms of overall economic losses in 2013 was the
41 flooding in southern and eastern Germany and neighbouring states in May and June with estimated
42 damages of \$15.2 billion. In June of the same year, flooding in India cost 5000 lives, with a further
43 2 million affected (MunichRe, 2014; EM-DAT).

44 The Global Assessment Report (UNISDR, 2011) states that the proportion of world population
45 living in flood-prone river basins increased by 114 percent over four decades from 1970 to 2010.
46 Additionally, while economic losses due to [river](#) floods have increased over the last 50 years, the
47 number of casualties has decreased. The reduction in loss of life has been associated with the
48 integration of early warning systems with emergency preparedness and planning at local and
49 national levels (Golnaraghi et al., 2009, Kundzewicz et al., 2012).

50 Global early warning systems are needed to improve international disaster management. These
51 systems can be used for both early forecasting, for better preparedness, and early detection, and
52 for an effective response and crisis management. Their necessity was emphasized in 2005, and
53 since then, it has been a key element of international initiatives such as the “Hyogo Framework
54 for Action 2005-2015” and, on a continental level, the European Commission Flood Action
55 Programme. After the 2002 floods on the Elbe and Danube rivers, the Commission supported the
56 development of the European Flood Awareness System (EFAS) (Bartholmes et al., 2009; Thielen
57 et al., 2009) by the Joint Research Centre to increase preparedness for riverine floods across
58 Europe. Currently, a number of organisations are involved in rapid mapping activities after major

59 (flood) disasters such as UNOSAT (2013), GDACS (2013), [“Space and Major Disasters” \(Disaster](#)
60 [Charter, 2014\)](#), [the Committee on Earth Observation Satellites \(CEOS\) Flood Pilot](#)-and the online
61 Dartmouth Flood Observatory (<http://floodobservatory.colorado.edu/>). In Europe, Copernicus is
62 the Earth Observation Programme which actively supports the use of satellite technology in
63 disaster management and early warning systems for improved emergency management.

64 Flood warning systems typically rely on forecasts from national meteorological services and in
65 situ observations from hydrological gauging stations. However, this capacity is not equally
66 developed across the globe, and is highly limited in flood-prone, developing countries. Ground
67 based hydro-meteorological observations are often either scarce or, in cases of transboundary
68 rivers, data sharing among the riparian nations can be limited or absent. Therefore, satellite
69 monitoring systems and global flood forecasting systems are a needed alternative source of
70 information for national flood authorities not in the position to build up an adequate measuring
71 network and early warning system. In recent years, there has been a notable development in the
72 monitoring of floods using satellite remote sensing and meteorological and hydrological modelling
73 (Schumann, et al. 2009).

74 A variety of satellite-based monitoring systems measure characteristics of the Earth’s surface,
75 including terrestrial surface water, over large areas on a regular basis (van Westen, 2013). Such
76 remote sensing is based on surface electromagnetic reflectance or radiance in the optical, infrared
77 and microwave bands. Some key advantages of microwave sensors is that they provide near-daily
78 basis global coverage and, at selected frequencies, relatively little interference from cloud cover.
79 Two presently-operating microwave remote sensors with near-global coverage are the Tropical
80 Rainfall Measuring Mission¹ (TRMM) operational from 1998 to present and the Advanced
81 Microwave Scanning Radiometer for Earth Observation System² (AMSR-E) which was active
82 from June 2002 to October 2011, followed by AMSR2 which was launched in May 2012 and is
83 onboard the Japanese satellite GCOM-W1³, and from which, brightness temperature data are being
84 distributed from January 2013 onwards. For future work, the European Space Agency (ESA) and
85 NASA have other missions to put similar instruments in orbit, capturing passive microwave energy

¹ <http://trmm.gsfc.nasa.gov>

² http://aqua.nasa.gov/about/instrument_amsr.php

³ http://suzaku.eorc.jaxa.jp/GCOM_W/w_amsr2/whats_amsr2.html

86 at 36.5 GHz, such as ESA's Sentinel-3 satellites (planned launch in 2015 and 2016) and NASA's
87 Global Precipitation Mission (GPM) (launched in February 2014) to replace TRMM.

88 Using AMSR-E data initially, De Groeve et al. (2006) implemented a method for detecting major
89 floods on a global scale, based on the surface water extent measured using passive microwave
90 sensing. Also, Brakenridge et al. (2005, 2007) demonstrated that orbital remote sensing can be
91 used to monitor river discharge changes. However, as underlined by Brakenridge et al. (2012,
92 2013), extracting the microwave signal and converting it into discharge measurements is not
93 straight-forward and depends on factors such as sensor calibration characteristics and perturbation
94 of the signal by land surface changes. These changes can be found for example in irrigated
95 agricultural zones and in areas where rivers flow along forested floodplains (Brakenridge et al.,
96 2013). As rivers discharge increases, river level (stage), river width, and river flow velocity all
97 increase as well, and the challenge is to measure one or more of these accurately enough to provide
98 a reliable discharge estimator, and compare against a background of other surface changes that
99 may affect what is measured from orbit.

100 There remains also the need to convert such discharge estimators to actual discharge units. Using
101 ground discharge data or climate-driven runoff models for calibration and validation, methods to
102 convert the remote sensing signal to river discharge have been previously tested at particular
103 stations with output from the Global Flood Detection System (GFDS,
104 <http://www.gdacs.org/floodddetection/>) and by different investigators (Brakenridge et al. 2007,
105 Brakenridge et al. 2012, Khan et al. 2012, Kugler and De Groeve, 2007, Moffitt et al. 2011, Hirpa
106 et al., 2013, Zhang et al. 2013). Yet the results are from different approaches and not easily
107 comparable, making an assessment of the potential performance on a global scale difficult.
108 Furthermore, definite conclusions about the influence of various environmental factors on the
109 signal performance have not been reached. Therefore, in this study, a rigorous broad assessment
110 of the method is undertaken with a systematic evaluation of the relationship between skills
111 obtained between ground- and satellite-based discharges, and the local characteristics of the
112 stations. Specifically this study addresses mean observed discharges, river widths, land cover
113 types, leaf area indices, climatic regions, and flood hazard maps, and the presence or absence of
114 large floodplains, wetlands, river ice and hydraulic control infrastructure.

115 Our goal is to assess the potentials and limitations of the satellite-based surface water extent signal
116 data for river discharge measurements with a large number of stations. Moreover, the relationship
117 between ground and satellite sets of discharge measurements and the local surface characteristics
118 is examined in order to provide guidelines for selection of observation sites. For this purpose,
119 river catchments located in a range of different climatic and land cover types were selected in
120 Africa, Asia, Europe, North America and South America. The remainder of the paper is structured
121 as follows: section 2 presents the study regions and data, section 3 describes the analysis
122 methodologies, and the results are discussed in section 4.

123

124 **2 Study regions and data**

125

126 **2.1. Study Regions and in situ discharge data**

127 Figure 1 shows the study basins and in situ discharge locations. The selected stations are all located
128 near major rivers of the world (Global Runoff Data Centre, 2007). The continental distribution and
129 the upstream catchment area of the stations are summarized in Table 1. We selected the locations
130 to be representative of a broad variety of local conditions: they belong to nine different main land
131 cover classes (aggregated from GlobCover, 2009) and five main types of climate (Peel et al., 2007).
132 The characteristics are listed in Table 2.

133 For Africa, Asia, Europe, North America and South America, daily in situ discharge values were
134 used from the Global Runoff Data Centre (GRDC) database. In addition, for the South African
135 stations, the discharge data were provided by the South African Water Affairs (DWA,
136 <http://www.dwa.gov.za/>). The selected stations for all these continents include daily data between
137 1998 and 2010, however not all stations have continuous data during this time period. From 1998,
138 the length of the time series was required to be above six years. The longest time series available
139 was of 13 years, with a median value of 8.5 years. In situ discharge information may itself be
140 affected by large and variable uncertainty, mostly on the measurement of the cross-sectional area
141 of the channel and mean flow velocity at the gauge or control site (Pelletier, 1988). Although
142 generally unknown, these values are typically between the 5-20% at the 95% confidence levels as
143 highlighted in studies such as [Hirsch and Costa \(2004\)](#), [Di Baldassarre and Montanari, \(2009\)](#),
144 [Le Coz et al. \(2014\)](#), and [Tominsk \(2014\)](#). However, the error uncertainty of river discharge is even

145 higher during floods events when the stage-discharge relationship, the so-called rating curve, is
146 used. As evaluated by Pappenberger et al. (2006), the analysis of rating curve uncertainties leads
147 to an uncertainty of the input of 18–25% at peak discharge. Di Baldassarre and Montanari (2009)
148 showed that the total rating curve errors increase, when the river discharge increases and varies
149 from 1.8% to 38.4% with a mean value of 21.2%. For the purposes here, these data are, however,
150 regarded as “ground truth”. We acknowledge the possible errors, however, and note that, for some
151 river reaches, satellite-based methods may actually track discharge changes more accurately than
152 ground-based measurements using stage; the extent to which this is true needs to be fully
153 investigates however.

154 (INSERT FIG 1 HERE)

155 (INSERT TABLE 1 HERE)

156 (INSERT TABLE 2 HERE)

157 **2.2. Satellite-derived data**

158 The Global Flood Detection System (GFDS) produces near real time maps and alerts for major
159 floods using satellite-based passive microwave observations of surface water extent and
160 floodplains. It is developed and maintained at the European Commission Joint Research Centre
161 (JRC) in collaboration with the Dartmouth Flood Observatory (DFO). The surface water extent
162 detection methodology using satellite-based microwave data is explained in Brakenridge et al.
163 (2007) and Kugler and De Groeve (2007). Here, only the basic principles are recalled.

164 At each pixel, the method uses the difference in brightness temperature, at a frequency of 36.5
165 GHz, between water and land surface to detect the proportion of within-pixel water and land. The
166 retrieved brightness temperature data are first gridded into a product with a pixel size of (near the
167 equator) 10 x 10 km (0.09 degree x 0.09 degree), and the system provides a daily output. For our
168 work, the merged TRMM/AMRS-E product was used
169 (<http://www.gdacs.org/flooddetection/download.aspx>); the gridded data are being provided in the
170 GCS WGS 1984 projection. For our period of study, 1998-2010, the merged data product was
171 employed for the time period of its availability (June 2002-2010), whereas stand-alone TRMM
172 data was used for the remaining time period (1998 to June 2002) and available latitudes. Note that

173 from 2013 the system is providing the merged product TRMM/AMSR2, however this period is
174 out of our scope.

175 In the GFDS system, the microwave signal (s) is defined as the ratio between the measurement
176 over wet pixel (M) and the measurement over a 7 pixel x 7 pixel array of background calibration
177 (C) pixel, [known as the M/C ratio](#)-(Brakenridge et al. 2012, De Groeve, 2010). Better discharge
178 signal values may be achieved when the measurement pixel is centred over a river reach and no
179 hydraulic structures are present (Moffitt, et al., 2011). However, this is sometimes difficult to
180 achieve due to the desired co-location with gauging stations (Brakenridge et al. 2012) or because
181 the potential measurement pixels within the raster are fixed, geographically.

182

183 **2.3. Other important datasets and maps**

184 The quality of the microwave signal detected by the satellite sensors can be influenced by local
185 ground conditions including extreme rainfall, snow/ice, land cover/use and topography
186 (Brakenridge et al., 2012). For example, forest is a type of land cover which influences the
187 microwave emission properties due to the biometric features of vegetation such as crown water
188 content and shape and size of leaves (Chukhlantsev, 2006). In this study, the effects of the local
189 ground conditions on the performance of the satellite signal were analysed as a function of the
190 following factors:

191 - **River width:** channel width from Yamazaki et al. (2014), estimation based on SRTM
192 Water Body Database and the HydroSHEDS flow direction map and for which the map
193 was upscaled from 0.025 to 0.1 degree, taking the mean of the river grid values in the 4 x
194 4 area.

195 - **Mean observed discharge:** For each station, a mean discharge value for the study period
196 was calculated from daily ground data (mainly from the GRDC dataset).

197 - **Upstream catchment area** (GRDC 2007) data: The GRDC river network was used to
198 visually select those stations located close to the “main rivers” classified by GRDC, and to
199 use the values of the upstream catchment area for each station. Note that upstream
200 catchment area values are missing from all South African stations from DWA data
201 provider.

- 202 - **Presence of Floodplains, Flooded Forest and Wetlands:** This was obtained from the
203 Global Lakes and Wetlands Database Level 3, a global raster map at 30-second resolution
204 which comprises lakes, reservoirs, rivers and different wetland types (Lehner and Doll,
205 2004).
- 206 - **Flood extent:** We used the fractional coverage of potential flooding of 25 km by 25 km
207 cells for a 100 year return period from the Global Flood Hazard Map derived using a model
208 grid (HTESSEL+CaMa-Flood) (Pappenberger et al. 2012).
- 209 - **Land cover:** We used land cover data from the Global Land Cover 2009 (GlobCover 2009)
210 (ESA and UCLouvain 2010). The 19 labels were aggregated into 8 types of land cover
211 depending on the vegetation type and density to synthesize the outputs (see Appendix Table
212 A 1). Further visual category checking was performed using GoogleMaps display for the
213 sites, and where necessary, land cover classes changed accordingly. An additional category
214 was added, for sparse vegetation areas where crops are grown along or near the river
215 channels.
- 216 - **Leaf Area Index:** A global reprocessed Leaf Area Index (LAI) from SPOT-VGT is
217 available for a period of 1999- 2007 (http://wdc.dlr.de/data_products/SURFACE/LAI/).
218 This LAI product is a global dataset of 36 ten-day composites at a spatial resolution of the
219 CYCLOPES products (1 km). For our analysis, a modified version of this product was
220 used, which was upscaled to a spatial resolution of 10 km.
- 221 - **Climatic areas:** We used the Köppen-Geiger climate map of the world (Peel et al. 2007)
222 to distinguish the main climate areas: tropical, arid, temperate, cold and polar (see Table
223 2).
- 224 - **Presence of river ice:** Through the signal, the presence of river ice cover can also be
225 detected in cold land regions. The Circum-Arctic Map of Permafrost and Ground-Ice
226 Conditions (Brown et al., 2002) map was used here. Examples of these rivers are Yukon
227 and Mackenzie in North America and Lena River in Russia. As is the case on the ground,
228 discharge under ice cover is left largely unmeasured as both water area and stage no longer
229 are responsive to discharge variation.

230 - **Dam location:** Hydraulic structures can disrupt the natural flow of water, and therefore
231 may alter the expected performance of the satellite signal on that location. For this analysis
232 the Global Reservoir and Dam (GRanD) (Lehner et al., 2008) dataset was used.

233

234 **3 Methodology**

235

236 **3.1. Satellite signal extraction**

237 In total, 398 locations for satellite-based measurement were selected which overlap spatially and
238 temporally with available in situ stations providing daily measurements. Since satellites never pass
239 directly over the same track at exactly the same time, the operational GFDS applies a four day
240 forward-running mean to systematically calculate ~~the M/C~~ signals; this also commonly fills
241 between any missing days (Kugler and De Groeve, 2007). Furthermore, for each observation site,
242 on the GFDS system the signal is calculated as the average signal of all measurement pixels under
243 observation for each location (which can be one or more pixels) (GDACS, 2014). Thus, in some
244 cases, even a 10 km pixel is not large enough as a measurement site, and would entirely saturate
245 with water during flooding. ~~A~~, an array of measurement pixels is instead used. In this analysis, we
246 used the signal values from the single pixels which contain the ground station, as well as a multiple
247 pixels selection. This includes, for each location, the pixel itself and also the three nearest
248 neighbours of the 10 x 10 km grid. In case of multiple pixels, the signal value was calculated for
249 the spatial median, average and maxima. Similar results were obtained globally when comparing
250 the extracted signals (single or multiple pixels) with the in situ discharge observations. Therefore,
251 we used the temporal and spatial averaging on the multiple pixel array as in the operational GFDS.
252 For each site, a visual check with Google maps was carried out to assure that the largest river
253 section was included within the finalized measurement sites (see Figure 2).

254 (INSERT FIGURE 2 HERE)

255

256 **3.2. Satellite signal calibration and validation**

257 For those co-located ground stations and satellite measurement sites where both sets of data (signal
258 and in situ discharge) were above six years in length, calibration and validation was performed

259 using the ground information as reference. Several stations, mainly in North America, located
260 close to man-made infrastructures such as weirs and generating stations were excluded from this
261 analysis due to the rapidly changing behaviour of the in situ observed discharge. Also, in a satellite-
262 based approach to measure river discharge, the local river characteristics and floodplain channel
263 geometry control the accuracy of rating curves as is the case for gauging stations on the ground
264 (Brakenridge et al., 2012, Khan et al., 2012 and Moffitt et al. 2011). Thus we expect some
265 measurement sites to exhibit a more robust response to discharge changes, and a higher signal to
266 noise ratio, than others.

267 It has been acknowledged that for large rivers, using the daily GFDS signal as a floodplain flow
268 surface area indicator of discharge might result in a few days lag when comparing with ground-
269 based discharge (Brakenridge, 2013). Thus, stage may immediately rise at a gauging station as a
270 flood wave approaches, but flow expansion out into the floodplain requires some increment of
271 time. This time lag may introduce error into the scatterplots used to calculate the rating equations,
272 and therefore lower skill scores obtained when analysing both datasets. In addition, in previous
273 studies (Khan et al. 2012, Zhang et al. 2013), it was observed that, in some cases, an overestimation
274 of satellite measured discharge existed during low flow periods when using a single rating equation
275 for the full period to calibrate signal into discharge units. For this reason, we decided to use a rating
276 equation for each month individually, and grouping daily into monthly data. In this case the time
277 series data for a fixed month can be treated as stationary and the derived daily discharge values
278 adjusted better also during low flow periods.

279 To calibrate satellite signal into discharge measurements, the first five years of data were used for
280 both satellite signal and ground discharge for each location. Regression equations were obtained
281 using monthly means from daily values and with which GFDS measured discharge was derived.

$$282 \quad Q_{\text{GFDSmeasured of X month}} = a_{\text{month}} + b_{\text{month}} * \text{signal} \quad (1)$$

283 For the sake of simplicity, for this paper, the equations were restrict to linear equations. However,
284 as the relation is purely empirical, we leave for follow on-work more research on flexible way to
285 fit these relations. Note that fitting straight lines to curves will reduce goodness of fit and predictive
286 accuracy.

287 The validation of the satellite derived daily discharge data was carried out with daily in situ data
288 on a two-year period, and skills scores were calculated to quantify the agreement between both
289 satellite and ground measured discharge. We are aware of the limited number of years (data) with
290 available time series for both variables, which might influence the robustness of the calibration.
291 In some cases there were longer time series available, but to standardised the analysis for all the
292 stations we used five years (1998-2002 or 2003-2008 for Northern stations with AMSR-E signal)
293 and the following two years for validation purposes (2003-2004 and 2009-2010 respectively). Note
294 that for 36 out of the 322 stations available data length was between six years and three months to
295 almost seven years. Validation was still carried out for the same period, but the data used for
296 calibration was slightly reduced. As an example, Figure 3a presents the scatterplot for the month
297 of March for the Senanga Station (Long 23.25, Lat. -16.116) in the Zambezi River (Africa) with
298 mean values derived from the period 1998 to 2002. For the same location, Figure 3b shows the in
299 situ observed and the GFDS measured discharge derived from the GFDS signal for the period
300 2003-2004.

301 (INSERT FIGURE 3 HERE)

302

303 **3.3. Skill scores**

304 The initial analysis of the correlation of the remote sensing signal to in situ discharge was assessed
305 for each station and site pair through the Pearson correlation coefficient (R). For the validation,
306 the performance of the satellite-measured discharge was also assessed using the Nash-Sutcliffe
307 Efficiency (NSE) statistic in addition to the R skill score.

308 One of the advantages of the R coefficient is its independence on the units of measurement, which
309 permits the comparison of dimensionless GFDS signal data. A small value indicates a weak or
310 non-linear relationship between the satellite signal and discharge. For this study, we grouped the
311 computed R values into three ranges as follows: <0.3 , $[0.3-0.7]$, and >0.7 .

312 Nash-Sutcliffe efficiency (NSE) (Nash and Sutcliffe, 1970) is typically used to assess the
313 predictive power of hydrological models and was here calculated to describe the accuracy of
314 satellite-derived discharge in comparison to gauge-observed discharge values. Higher values of
315 the Nash-Sutcliffe statistic should indicate more correlated results, without other factors taken into
316 account, such as autocorrelation (Brakenridge et al., 2012). However, the degree of correlation of

317 these variables does not verify the discharge magnitudes (Brakenridge et al., 2013). A NSE value
318 of 1 corresponds to a perfect match of modelled to the observed data whereas NSE = 0 indicates
319 that the model predictions are as accurate as the mean of the observed data. ~~Thus here model~~
320 ~~simulations are judged as “satisfactory” if NSE > 0.50 (Moriassi et al., 2007).~~ The resulting scores
321 will be classified as in Zaraj, et al. (2013): < 0, [0.2-0.5], [0.5-0.75], and > 0.75.

322

323 **3.4 Factors affecting the satellite signal**

324 Understanding the influence of local factors on the accuracy of the satellite flood detection is
325 critical for practical use of the remotely sensed signal. We analysed the accuracy effects of river
326 width, mean daily discharge, upstream catchment area, presence of large floodplain, flooded forest
327 and wetlands, the potential flood extent, land cover type, Leaf Area Index (LAI), climatic areas,
328 presence of river ice and hydraulic structures. To assess their influence, the fractional coverage
329 over the measurement site was retrieved for variables with spatial coverage.

330 First, we use the skill scores (R and NSE) obtained from a simple analysis for each individual
331 factor or variable. Second, we seek to understand which of the surface variables have the greatest
332 importance in determining sites with a good or poor performance. For this purpose, we use a
333 decision tree technique called Random Forest (RF). Among other features, this allows ranking of
334 the relative importance of each variable. The technique is described by Breiman (2001) and
335 implemented in R by Liaw and Wiener (2002), where the reader is referred for a more detailed
336 explanation. As a summary of the Random Forest algorithm, *n*tree bootstrap samples are randomly
337 selected from the data set, a different subset is used for each bootstrap and for each sample a tree
338 is grown, obtaining *n*tree trees. Random Forest is called an ensemble method because it applies
339 the method for a number of decision trees, in this case 500, in order to improve the classification
340 rate. Some stations are left out of the sample (out-of-bag) and used to gain an internal unbiased
341 estimate of the generalisation error (oob errors) and to obtain estimates of the importance of the
342 variables (Breiman, 2001). These values are averaged over the *n*tree trees. For the variables
343 classification, the node impurity is measured by the Gini index. [Gini's mean difference was first](#)
344 [introduced by Corrado Gini in 1912 as an alternative measure of variability and the parameters](#)
345 [derived from it, such as the Gini index, also referred to as the concentration ratio \(Yitzhaki and](#)
346 [Schechtman,2013\). The Gini index is mostly popular in economics, however it is also used in other](#)

347 [areas, such as building decision trees in statistics to measure the purity of possible child nodes,](#)
348 [and it has been compared with other equality measures \(Gonzalez, L., et al. 2010\). This index is one](#)
349 [of the most frequently used measures of heterogeneity for selecting the best splitting variable](#)
350 [\(Sandri and Zuccolotto, 2008\).](#) The variables with higher decrease in Gini values (lower Gini) are
351 those with higher importance on the classification analysis.

352 Although for “black-box models” such as Random Forest the information is hidden inside the
353 model structure, the prediction power is high (Palczewska et al., 2013). This method is relatively
354 robust given outliers and noise because it uses randomly chosen subsets of variables at each split
355 of each tree (Breiman, 2001; Chan et al., 2008). To further increase robustness, Strobl et al. (2009)
356 states that results from the random forest and conditional variable importance should always be
357 tested by doing multiple random forest runs using different seeds and sufficiently large ntree values
358 to obtain robust and stable results.

359 The quality index chosen to rank variable importance and classify good or poor locations, in the
360 Random Forest analysis, was the Nash-Sutcliffe Efficiency (NSE) score. A threshold of NSE=0
361 splits the data into two groups, obtaining about 50% of the data above (true or good predictive)
362 and below (false or poor predictive) that value of NSE. The results presented here are the average
363 of 200 runs. Furthermore, four different training sets were used by a random 70%/75%/80%/90%
364 of the stations and [wereas](#) validated with the remaining 30%/25%/20%/10% of stations,
365 respectively.

366

367 **4 Results and discussion**

368

369 As a first step we analysed the relationship between the satellite signal and the in situ observed
370 discharge to have an initial understanding of the performance between the two datasets (Section
371 4.1). Then we calibrated the satellite signal with in situ discharge data. With the regression
372 equations obtained, we calculated [satellites](#) discharge ~~satellites~~ measurements. A two-year
373 validation period was carried out for each station using the skill scores as described in Section 3.3
374 (Section 4.2). This was followed by an assessment for how different variables contribute in a
375 positive or negative way to the overall skill (Section 4.3). Variables included in the analysis are

376 daily mean river discharge, river width, upstream catchment area, potential flood hazard area, land
377 cover, leaf area index, climatic zones, presence of large floodplains, flooded forest and wetlands,
378 river ice and hydrologic structure. Finally, the relative importance of all variables in comparison
379 to each other has been assessed (Section 4.4).

380 Before analysing the validation results, it is important to highlight two possible different sources
381 of error which might influence the outputs. Firstly, the signal ~~may be to noise ratio noisy in general~~
382 ~~might be low~~ for a site or have intermittent instrument noise occasionally producing intermittent
383 positive spikes in discharge. ~~large noise values (instrument noise) coming from the raw signal data.~~
384 Secondly, the rating curve may be offset, which will result in a consistent bias on the discharge
385 values for that location even though the time series are strongly correlated.

386

387 **4.1. Correlation of raw satellite data vs. gauge observations**

388 The first step was to look at the “raw” correlation between daily ground station-measured water
389 discharge and the satellite signal and to calculate the empirical linear relation between these two
390 variables for each site. The full time series, including low flows, were used for the calculation and
391 executed for 398 stations. Figure 4 shows the R skills obtained. 169 ~~stations~~ out of 398 sites have
392 an $R > 0.3$ and 42 of them have $R > 0.5$. Perhaps, correlations might have been higher if regression
393 would have not been restricted to linear equations (Brakenridge et al., 2007, 2012).

394 (INSERT FIGURE 4 HERE)

395

396 **4.2. Satellite signal calibration, validation and evaluation through skill scores**

397 For the stations with over six years of contemporary data for both in situ discharge and satellite
398 signal, we obtained regression equations for each month of the year and station using the first five
399 years of data. Next, using these equations we carry out a calibration of the daily signal into
400 discharge units. Afterwards, the validation of the GFDS measured discharge was implemented for
401 the following two years. In some regions such as Northern Asia, the lack of available recent long
402 time series (after 2002) meant that the number of stations available for calibrating the satellite into
403 discharge measurements was reduced. Stations where the number of years matching observed
404 discharge and satellite signal was shorter rd than six years were excluded from the validation

405 exercise despite performing well. Finally, out of 398 a total of 332 stations remained for calibration
406 and validation.

407 Figure 5 shows that for NSE score, 154 out of 332 stations are larger than 0; 13 located in Africa,
408 77 in North America, 62 in South America, 1 in Asia and 1 in Europe. Nevertheless, it needs to be
409 noted that in arid regions, results calculated with the skill scores such as NSE are penalised, by
410 low average discharge compared to high flow conditions. If instead of using all the available time
411 series, a “dry stream” threshold would have been applied, the scores obtained for these sites could
412 have been higher when analysing the remaining dataset period where flow is present.

413 (INSERT FIGURE 5 HERE)

414

415 **4.3. Analysis of the factors affecting the satellite signal**

416

417 **4.3.1. River width and presence of floodplain and wetlands.**

418 As a first step to analyse the potential relationship between the individual local characteristics and
419 the performance of the locations in global terms, we study the R score of the validation for the 322
420 stations in relation with the maximum river width value at each location (Figure 6a). Results
421 indicate that locations with a river width higher than 1 km are more likely to score an R larger than
422 0.3. In fact, the mean R score is 0.60. Where 26 out of 64 (~41%) have R> 0.75. However, there
423 is a number of stations with lower river width that also obtained high scores. As the retrieval of
424 the satellite signal also depends on the floodplain geometry. As soon as the river floods and water
425 goes over-bank, the proportion of water in the wet pixel greatly increases. So the score should be
426 also high for small rivers with a proportionally big floodplain. Figure 6b shows the R scores by
427 locations where the majority of the area belongs to floodplain, flooded forest and wetlands
428 category or, their absence. In our study, higher median scores were obtained for those located in
429 large freshwater marsh and floodplains, followed by those on swamps and flooded forest. These
430 results give a first indication on the characteristics of the locations with better performance.

431 (INSERT FIGURE 6 HERE)

432 **4.3.2. River discharge and potential flooding**

433 Flooding is determined by the discharge as well as the potential flood hazard. Figure 7a shows that
434 84 out of 95 stations with R<0.3, also have mean discharge values lower than 500 m³s⁻¹ (Log10

435 (500) ≈ 2.7), of which 55 stations in fact had a mean discharge lower than $200 \text{ m}^3\text{s}^{-1}$. These stations
436 are mainly located in South Africa, and in some areas of North America. Therefore, ~~it~~ can be
437 concluded that the mean discharge can be considered a key variable that determines the
438 appropriateness of locations for which satellite discharges can be derived: As 77% of the stations
439 with $Q < 500 \text{ m}^3/\text{s}$, have $R < 0.3$, while 91.5% of the stations with $Q > 500 \text{ m}^3/\text{s}$ have $R > 0.3$, locations
440 with discharge of less than $500 \text{ m}^3\text{s}^{-1}$ might not provide reliable results for a global satellite-based
441 monitoring system. Alternatively, non-permanent rivers and streams exhibiting only seasonal or
442 ephemeral flow (typical for dry regions) may require a different monitoring approach, wherein a
443 “dry” threshold is established for the signal data.

444 After excluding the global stations with low skill score due to low flows and studying the
445 remaining stations, we can better understand the performance of the system in relation to other
446 local characteristics. Figure 7b shows for each location the relationship between the validation R
447 and the percentage of area in each pixel covered by potential flooding during a 100 year return
448 period flood event, obtained with the model grid (HTESSEL+CaMa-Flood) (downscaled from a
449 $25 \times 25 \text{ km}$ pixel, Pappenberger et al., 2012). 100 means totally flooded across its area, 50 means
450 50 % of the area within the cells is flooded, and 0 means that the area is not flooded. Although
451 there is not a clear trend for all the points, result indicate that locations with a percentage of
452 potential flooding larger than 40%, are expected to score an R larger than 0.3.

453 (INSERT FIGURE 7 HERE)

454 **4.3.3. Land cover types and climatic areas**

455 Figure 8 presents a global evaluation of the R score obtained during the validation and its
456 classification by the land cover type of the stations. The bare land cover category was excluded
457 from this study as only one of the selected locations belong to that class. Looking at the median of
458 the boxplot (see Figure 8), we found that some of the locations with higher density of vegetation
459 such as those located on “closed forest” and “mosaic with predominant vegetation” (included
460 forest, scrublands and grasslands) obtained lower median scores values. In contrast, the locations
461 with lower vegetation density such as “sparse vegetation”, “mosaics with predominant
462 cropland/grasslands”, “open forest” and “closed to open forest” land cover types obtained larger
463 median R scores, around 0.6-0.8. Similar results can be observed when looking at the interquartile
464 range or spread of the boxplots: “closed to open forest” and “mosaics with predominant

465 cropland/grasslands” obtained better results. Meanwhile, “closed forest” and “mosaic with
466 predominant vegetation” had lower scores. In addition, those sites with a combination of sparse
467 vegetation and crops growing near the river channel had a lower median value when comparing
468 with those on sparse or mosaic crops land cover. ~~To~~Note that the sites with “sparse with crops”
469 are located in arid climatic areas, whereas most of the “sparse” are in cold or polar regions,
470 therefore run by different processes. In addition, sites with a majority of artificial/urban land cover
471 (not shown) obtained a low median value of 0.267.

472 (INSERT FIGURE 8 HERE)

473

474 The relationship between locations by main Köppen-Geiger climatic areas (Peel et al. 2007) and
475 R score obtained is shown in Figure 9. Globally the tropical regions (Africa and South America)
476 obtained the highest median scores ($R \approx 0.8$), followed by cold regions ($R \approx 0.6$). Lower median
477 score values ($R \approx 0.3$) were obtained for arid and temperate regions. It is important to clarify that
478 these results are not only due to direct climate characteristics but also for example due to the
479 characteristics of the rivers on those areas. In the case of the arid regions, it is mainly related with
480 reduce daily average discharges, a characteristic of many of these stations. Note that polar climate
481 was excluded from this evaluation as only three locations belong to that class.

482 (INSERT FIGURE 9 HERE)

483 **4.3.4. Leaf Area Index (LAI)**

484 Leaf Area Index (LAI) values typically range from 0 for bare ground to 6 or above for a dense
485 forest, however CYCLOPES underestimates over dense vegetation (forest) (Zhu et al., 2013).
486 Therefore, for this product LAI range is limited to [0-4], as seen in our analysis. Despite this,
487 CYCLOPES is the most similar product to LAI references map (*Ibid.*). According to the study
488 carry out by Zhu et al. (2013) monthly CYCLOPES LAI values for the period 1999 to 2007 by
489 four main groups of vegetation are predominantly as follows: bare ground [0], forest [0-3.5], other
490 woody vegetation [0-1.5], herbaceous vegetation [0-2], and cropland/natural vegetation mosaics
491 [0-3]. ~~Where~~The highest annual mean LAI values are obtained by evergreen broadleaf forest
492 (3.16), included in our “closed to open forest” class.

493 We decided to study the relationship between the mean Leaf Area Index and the skill obtained in
494 the validation for each location, also looking at complementary variables such as the land cover
495 and the geographical region which the stations belong to. Figure 10 shows that locations with a
496 mean [LAI > 2] predominantly have a “closed to open forest” type in South America (31 stations)
497 of which 29 have an R score higher than 0.6. For [LAI > 2] there is also 12 North American
498 locations with “closed forest” land cover but in general with poorer scores for those locations.
499 Additionally, 18 stations with mosaic vegetation from North and South America obtained [LAI >
500 2] and 16 out of them, a [R>0.6]. For [LAI < 2], both the land cover and geographical locations
501 are distributed along the scatterplots, from poor to high correlations.

502 (INSERT FIGURE 10 HERE)

503

504 **4.3.5. River ice**

505 Figure 11a shows the scores obtained for the locations with presence or ~~not~~-absence of river ice,
506 including a range from continuous to sporadic (Brown et al., 2002). It can be seen that stations
507 located in areas with river ice tend to have a good correlation between in situ and satellite
508 measured discharge (based on 33 stations), as the system tends to capture well the annual spring
509 ice break-up and freezing as indicated in the study by Brakenridge et al.(2007) and Kugler (2012).
510 At these locations, once ice-covered there is no sensing capability from the system: which may
511 seem analogous to low flow conditions, and for which sites we obtained lower scores. However,
512 there is an important difference when analysing time series of signal between ice covered high
513 latitude river and all-year-around low flow rivers. When on the sites with river ice melting process
514 takes place, there is an increase of runoff happening and for many places the signal strongly
515 indicates this increased flow. On the other type of rivers, low flows is generally a characteristic for
516 ~~the~~-most of- the year and if the signal to noise is low, the signal retrieved is very noisy: one
517 motivation for setting a “dry” threshold for such sites.

518

519 **4.3.6. Hydraulic structures**

520 The correlation between satellite and discharge data depends on both variables. Typically it is
521 assumed that observed discharges are “ground truth”, however, when influenced by structures and
522 dams the ground discharge may not be well-monitored by flow area/flow width variation. For
523 example, when there is a major increase in river discharge but a flood is avoided by artificial

524 levees, we cannot expect that the satellite signal will accurately capture the flood hydrograph; as
525 well, downstream flooding may be attenuated by an upstream flood control dam and reservoir; so
526 that the gauge location is critical. Figure 11b shows the influence of the presence or absence of a
527 nearby dam using the Global Reservoir and Dam (GRanD) database (Lehner et al., 2008) or
528 visually identified hydraulic control infrastructure. Locations where the dam or other element was
529 present (48 stations) obtained lower median R score. Therefore, ideally, observation sites should
530 be located in areas without hydraulic control infrastructures.

531 (INSERT FIGURE 11 HERE)

532

533 **4.4. Variable importance**

534 Based on the individual analysis of the signal potential influence factors we found that to
535 understand the site performances, in some occasions multiple variables need to be analysed in a
536 simultaneous way. For example, [the general low scores obtained on the Eastern USA stations](#)
537 [might be due to a number of factors: ~64% of these stations have a mean discharge value lower](#)
538 [than 500 m³s⁻¹ and ~88% of the stations are located at river width lower than 1km. In addition,](#)
539 [~59% of the stations are located in wetlands areas. Another example, in this case](#) regarding the
540 exceptions of the low R and mean observed discharge higher than 500 m³s⁻¹, all the 11 locations
541 have a potential probability of flooding lower than 21%, the land cover of 10 out of 11 is forest, 5
542 of them located in wetlands and two of them have a nearby hydraulic structure. Despite exhibiting
543 a mean discharge greater than 500 m³s⁻¹, these other local characteristics may be the cause of the
544 poor performance. Therefore, we decided to use a classification decision tree technique (Random
545 Forest), which split the dataset at each node according to the value of one variable at a time (the
546 best split) from a selected set of variables to understand the importance of each variable. Random
547 Forest is called an ensemble method because it is performed for a number of decision trees, in this
548 case 500 trees, in order to improve the classification rate.

549 The result presented here is the rank of the importance of variables to classify a location with a
550 good or poor performance. These values are obtained as an output of the Random Forest analysis
551 and are, in addition, the average of 200 independent runs. As explained in section 3.4 the variable
552 importance based on the mean decrease in Gini index was calculated for the Nash-Sutcliffe
553 Efficiency (NSE) score obtained from the validation. We used a NSE=0 to distinguish the sites

554 with a good (above 0) from poor performance (below 0) and we also tested it with a threshold NSE
555 of 0.50.

556 Figure 12 presents the variable importance for the four test groups. Features which produced large
557 values of the “Mean Decrease in Gini” are ranked as more important than features which produced
558 small values. For our locations and data available the mean daily observed discharge has the
559 highest importance, followed by the climatic region, land cover / mean LAI and upstream
560 catchment area. Meanwhile, the presence of hydraulic structures (mainly dams) and of river ice
561 has the lowest importance to classify a location as good or poor performance. However, this does
562 not mean that it has no influence. Although discharge is correlated with upstream catchment area
563 and at some degree also leaf area index with land cover type, both were included in this case to
564 understand which variable might help us most to classify the sites.

565 Although, the effect of the correlations on these measures has been studied recently (see Archer
566 and Kimes (2008), Strobl et al. (2008), Nicodemus and Malley (2009), Nicodemus et al. (2010),
567 Nicodemus (2011), Auret and Aldrich (2011), Tolosi and Lengauer (2011), Grömping, U. (2009)
568 and Gregorutti et al. (2013)) there is no yet a consensus on the interpretation of the importance
569 measures when the predictors are correlated and on what is the effect of this correlation on the
570 importance measure.

571 In order to test the effect on the results when correlated variables were included in our analysis, an
572 independent Random Forest analysis was carried out (not shown in the paper) for the same
573 variables but excluding the river width and the presence of floodplains and wetlands variables.
574 Results also showed that the mean daily observed discharge had the highest importance and the
575 presence of hydraulic structures (mainly dams) and of river ice had the lowest importance to
576 classify a location as good or poor performance.

577 (INSERT FIGURE 12 HERE)

578

579 **5 Conclusions and future research**

580

581 In this article we presented an evaluation of the skill of the Global Flood Detection System to
582 measure river discharge from remote sensing signal. From the 322 stations validated the average

583 continental R skills are as follow: Africa 0.382, Asia 0.358, Europe 0.508, North America 0.451
584 and South America 0.694. Approximately 48% of these stations have an NSE score higher than
585 zero; 13 located in Africa, 77 in North America, 62 in South America, 1 in Asia and 1 in Europe.
586 Results showed that the majority of the stations that received low skills scores, were due to low
587 flow conditions. For example, 84 out of 95 stations with $R < 0.3$, have mean discharge values lower
588 than $500 \text{ m}^3\text{s}^{-1}$. These are located mainly in South Africa with 25 cases and North America with
589 53 cases, which penalised their average continental skills. Note that our focus was on factors
590 affecting the method, globally, and that these skill values do not directly indicate at-a-site
591 measurement accuracy (which could be improved, for example, by use of non-linear rating
592 equations and/or accommodation of any phase shift or timing differences in flow area- versus state-
593 based discharge monitoring).

594 In order to better understand the impact of the local conditions on the performance of the sites, we
595 looked first at specific factors individually. In general terms, higher skill scores were obtained for
596 location with one or more than one of the following characteristics: a river width higher than 1km,
597 a large floodplain area, in flooded forest, with a potential flooded area per pixel greater than 40%
598 during a 100 year return period flood event, a land cover type of sparse vegetation, croplands or
599 grasslands and closed to open and open forest, Leaf Area Index above 2, located in a tropical
600 climatic area and where no dams or hydraulic infrastructures are present. Also, out of our
601 locations, high latitude rivers with seasonal ice-cover tend to exhibit good performance.

602 Secondly, we performed a classification decision tree analysis, based on Random Forest, to obtain
603 the variable importance when classifying a site as good or poor. The output of this analysis showed
604 that mean observed discharge, climatic region, land cover and mean leaf area index (LAI) and
605 upstream catchment area and were the variables with higher importance, whereas river ice and
606 dam obtained the lowest importance. Both the individual and the combined classification analysis
607 of these local characteristics give us critical evidence of the relationship between the ground and
608 satellite discharge measurements and when it is expected to perform well. Furthermore, it provides
609 a guideline for future selection of measuring sites.

610 The locations with a very good performance will be selected for a potential future project where
611 satellite measure discharge could be calculated for longer periods and on a daily basis from the
612 remote sensing signal, analogous to the Dartmouth Flood Observatory method. This will represent

613 a major step forward in developing continental and global hydrological monitoring systems as
614 these data can fill the gaps where real time ground discharge measurements are not available (the
615 case at many locations globally). We found that some of the sites with good performance are
616 located within international river basins such as the Niger, Volta and Zambezi in Africa. In
617 addition, for the studied locations with good signal performance but rather short contemporary
618 time series with in situ observed discharge (such as the Siberian stations), the calibration of the
619 signal to obtain discharge measurements could be executed at any point when additional ground
620 data is available. This will also be beneficial for all stations including those with time series above
621 seven years long.

622 Zhang et al. (2013) recently demonstrated the potential of integrating satellite signal provided by
623 the Global Flood Detection System in improving flood forecasting. This first attempt of data
624 assimilation was carried out for a single station (Rundu, northern Namibia- included in this study)
625 with the conceptually simple Hydrological MODel (HyMOD). Hence, a prospective study with
626 the inclusion of all these stations for post-processing through data assimilation and error correction
627 of the stream-flow forecast in hydrological models could be done. For instance, for the pre-
628 operational Global Flood Awareness System (GloFAS) (Alfieri et al. 2012) and the African Flood
629 Forecasting System (AFFS) (Thiemig et al. 2014) in an analogous way as it is already being done
630 with ground gauge observed streamflow on the European Flood Awareness System (Bartholmes
631 et al., 2009; Thielen et al., 2009). Hence, work towards the integration of global flood detection
632 and forecasting systems such as GFDS and GloFAS, respectively, can provide a more
633 comprehensive information for decision makers.

634
635 *Acknowledgements.* We acknowledge the Global Runoff Data Centre and South African Water
636 Affairs for providing historic discharge measurements. Furthermore we would like to acknowledge
637 the team from the Joint Research Centre Crisis Monitoring and Response Technologies
638 (CRITECH) for support and access to the Global Flood Detection System signal historical data.
639 Also, Philippe Roudier, Simone Russo, Angel Udias and Feyera Hirpa are thanked for their
640 valuable input and methodology advice and Ad de Roo for PhD supervision and the editor and the
641 two reviewers are gratefully acknowledged for their valuable feedback. G. R. Brakenridge
642 acknowledges funding support from the NASA Hydrology Program.

References

- Alfieri, L., Burek, P., Dutra, E., Krzeminski, B., Muraro, D., Thielen, J. and Pappenberger, F.: GloFAS-global ensemble streamflow forecasting and flood early warning, *Hydrology and Earth System Sciences*, vol. 17, no. 3, pp. 1161-1175, 2013.
- [Archer, K. J. and Kimes, R. V. Empirical characterization of random forest variable importance measures. *Computational Statistics and Data Analysis*, 52:2249–2260, 2008. doi: 10.1016/j.csda.2007.08.015](#)
- [Auret, L. and Aldrich, C. Empirical comparison of tree ensemble variable importance measures. *Chemometrics and Intelligent Laboratory Systems*, 105:157–170, 2011. doi: 10.1016/j.chemolab.2010.12.004](#)
- Bartholmes, J.C., Thielen, J., Ramos, M.H. and Gentilini, S.: The European flood alert system EFAS - Part 2: Statistical skill assessment of probabilistic and deterministic operational forecasts. *Hydrology and Earth System Sciences*, 13(2): 141-153, 2009.
- Brakenridge, G. R., S. V. Nghiem, E. Anderson, and S. Chien, Space-based measurement of river runoff, *Eos Trans. AGU*, 86(19), 185–188, 2005, doi:10.1029/2005EO190001.
- Brakenridge, G.R., Nghiem, S.V., Anderson, E. & Mic, R.: Orbital microwave measurement of river discharge and ice status, *Water Resources Research*, vol. 43, no. 4, 2007, W04405, doi:10.1029/2006WR005238.
- Brakenridge, G.R., Cohen, S., Kettner, A.J., De Groeve, T., Nghiem, S.V., Syvitski, J.P.M. and Fekete, B.M.: Calibration of satellite measurements of river discharge using a global hydrology model, *Journal of Hydrology*, vol. 475, pp. 123-136, 2012.
- Brakenridge, G.R., De Groeve, T., Cohen, S., and Nghiem, S. V.: River Watch, Version 2: Satellite River Discharge and Runoff Measurements: Technical Summary, University of Colorado, Boulder, CO, USA, <http://floodobservatory.colorado.edu/SatelliteGaugingSites/technical.html>, last access: 1 December 2013.
- Breiman, L.: Random Forests. *Machine Learning*, 45, 5–32, 2001.
- Brown, J., O.J. Ferrians, Jr., J.A. Heginbottom, and E.S. Melnikov.: Circum-Arctic Map of Permafrost and Ground-Ice Conditions. Version 2. [Permafrost], Boulder, Colorado USA: National Snow and Ice Data Center, 2002.

1 [Committee on Earth Observation Satellites \(CEOS\) Flood Pilot, http://www.ceos.org/](http://www.ceos.org/), last access:
2 [1 September 2014.](#)

3 Chan, J.C.-. & Paelinckx, D.: Evaluation of Random Forest and Adaboost tree-based ensemble
4 classification and spectral band selection for ecotope mapping using airborne hyperspectral
5 imagery, *Remote Sensing of Environment*, vol. 112, no. 6, pp. 2999-3011, 2008.

6 Chukhlantsev, Alexander A.: Modeling of microwave emission from vegetation canopies,
7 *Microwave Radiometry of Vegetation Canopies*. Springer Netherlands, Chapter 6. pp. 147–
8 175, 2006.

9 [Di Baldassarre, G. and Montanari, A.: Uncertainty in river discharge observations: a quantitative
10 analysis, *Hydrol. Earth Syst. Sci.*, 13, 913-921, doi: 10.5194/hess-13-913-2009, 2009.](#)

11 De Groeve, T., Brakenridge, G. R., and Kugler., Z.: Near Real Time Flood Alerting for the Global
12 Disaster Alert and Coordination System, eds B. Van de Walle, P. Burghardt, and C.
13 Nieuwenhuis *Proceedings of the 4th International ISCRAM Conference*, 33-40, 2006.

14 De Groeve, T., and Riva, P.: Global Real-time Detection of Major Floods Using Passive
15 Microwave Remote Sensing, *Proceedings of the 33rd International Symposium on Remote
16 Sensing of Environment Stresa, Italy, May 2009.*

17 De Groeve, T.: Flood monitoring and mapping using passive microwave remote sensing in
18 Namibia, *Geomatics, Natural Hazards and Risk*, 1:1, 19-35, 2010.

19 Di Baldassarre, G. & Montanari, A.: Uncertainty in river discharge observations: A quantitative
20 analysis, *Hydrology and Earth System Sciences*, vol. 13, no. 6, pp. 913-921, 2009.

21 [Disaster Charter, 2013. Space and Major Disasters. http://www.disasterscharter.org/](http://www.disasterscharter.org/), last accessed 1
22 [September 2014.](#)

23 EM-DAT, The OFDA/CRED International Disaster Database, Université Catholique de Louvain,
24 Brussels, Belgium, <http://www.emdat.be>, last access: 1 December 2013.

25 Fekete, B.M., Vorosmarty, C.J., Grabs, W., 1999. Global, composite runoff fields based on
26 observed river discharge and simulated water balances, GRDC Report 22, Global Runoff
27 Data Center, Koblenz, Germany.

28 GDACS, Global Disaster Alert and Coordination System, Global Floods Detection System.
29 <http://www.gdacs.org/>, last accessed 1 December 2013.

30 Global Runoff Data Centre: Major River Basins of the World, 2007. 56068 Koblenz, Germany:
31 Federal Institute of Hydrology (BfG). <http://grdc.bafg.de/>, last accessed 20 January, 2013

- 1 Global Runoff Data Centre, The. River Discharge Time Series. 56068 Koblenz, Germany: Federal
2 Institute of Hydrology (BfG). <http://grdc.bafg.de/>, last accessed 20 January, 2013
- 3 [Golnaraghi M., J. Douris, J.-B. Migraine: Saving Lives Through Early Warning Systems and
4 Emergency Preparedness, Risk Wise, Tudor Rose, pp 137–141, 2009.](#)
- 5 [Gonzalez, L., Velasco Morente,F. , Gavilan Ruiz, J.M., Sanchez-Reyes Fernandez, J.M. The
6 Similarity between the Square of the Coeficient of Variation and the Gini Index of a
7 General Random Variable. Journal of Quantitative Methods for Economics and Business
8 Administration 10: 5–18.2010, ISSN 1886-516X.](#)
- 9 [Gregorutti,B., Michel, B., Saint-Pierre, P. Correlation and variable importance in random forests.
10 Cornell University Library, 2013. arXiv: 1310.5726 \[stat\]](#)
- 11 [Grömping, U. Variable Importance Assessment in Regression: Linear Regression versus Random
12 Forest. The American Statistician. 11/2009; 63:308-319, 2009. doi: 10.1198/tast.2009.08199](#)
- 13 Hirpa FA, Hopson TM, De Groeve T, Brakenridge GR, Gebremichael M, Restrepo PJ. Upstream
14 satellite remote sensing for river discharge forecasting: Application to major rivers in South
15 Asia. Remote Sens Environ, 131:140-51, 2013.
- 16 Hirsch R. M, Costa J. E.:U.S. Stream Flow Measurement and Data Dissemination Improve EOS,
17 Transactions, American Geophysical Union. Vol. 85, No. 20, 18 May 2004, 197-203 pp,
18 2004.
- 19 Khan, S.I., Hong, Y., Vergara, H.J., Gourley, J.J., Robert Brakenridge, G., De Groeve, T., Flamig,
20 Z.L., Policelli, F. & Yong, B.: Microwave satellite data for hydrologic modeling in ungauged
21 basins, IEEE Geoscience and Remote Sensing Letters, vol. 9, no. 4, pp. 663-667, 2012.
- 22 Kugler, Z., and De Groeve, T.: The Global Flood Detection System, Office for Official
23 Publications of the European Communities, Luxembourg, 2007.
- 24 Kugler, Z.: Remote sensing for natural hazard mitigation and climate change impact assessment,
25 Quaterly Journal of the Hungarian Meteorological Service. Vol.116, No.1, January-March
26 2012, pp.21-38, 2012.
- 27 [Kundzewicz, Z. W: Changes in Flood Risk in Europe, Wallingford: IAHS Press. 516 p. IAHS
28 special publication; 10, 2012. United Nations: Report of the United Nations Conference on
29 Sustainable, Development Rio de Janeiro, Brazil. 20–22 June 2012, A/CONF.216/16, 2012.](#)

- 1 Le Coz, J., Renard, B., Bonnifait, L., Branger, F. & Le Boursicaud, R.: Combining hydraulic
2 knowledge and uncertain gaugings in the estimation of hydrometric rating curves: A
3 Bayesian approach, *Journal of Hydrology*, vol. 509, pp. 573-587, 2014.
- 4 Lehner, B., and Döll, P.: Development and validation of a global database of lakes, reservoirs and
5 wetlands. *Journal of Hydrology* 296/1-4: 1-22, 2004. doi: 10.1016/j.jhydrol.2004.03.028
- 6 Lehner, B., Reidy Liermann, C., Revenga, C., Vörösmarty, C., Fekete, B., Crouzet, P., Döll, P.,
7 Endejan, M., Frenken, K., Magome, J., Nilsson, C., Robertson, J., Rödel, R., Sindorf, N.,
8 Wisser, D.: High resolution mapping of the world's reservoirs and dams for sustainable river
9 flow management, *Frontiers in Ecology and the Environment*. Source: GWSP Digital Water
10 Atlas. Map 81: GRanD Database (V1.0), 2008. Last access: 11/03/2014.
11 http://atlas.gwsp.org/index.php?option=com_content&task=view&id=209&Itemid=1
- 12 Liaw, A. and Wiener, M. Classification and Regression by randomForest. *R News* 2(3), 18—22,
13 2002.
- 14 Moffitt, C.B., F. Hossain, R.F. Adler, K.K. Yilmaz, and H.F. Pierce. Validation of a TRMM-Based
15 Global Flood Detection System in Bangladesh. *International Journal of Applied Earth
16 Observation and Geoinformation* Volume 13, Issue 2, April 2011, Pages 165-177, DOI:
17 10.1016/j.jag.2010.11.003.
- 18 MunichRe, Munich Reinsurance: January 2014 press release, 2014.Münchener
19 Rückversicherungs-Gesellschaft, Geo Risks Research, NatCatSERVICE
20 http://www.munichre.com/en/media_relations/press_releases/2014/2014_01_07_press_rele
21 [ase.aspx](http://www.munichre.com/en/media_relations/press_releases/2014/2014_01_07_press_rele) , last access 20 January 2014
- 22 Nash, J.E. & Sutcliffe, J.V.: River flow forecasting through conceptual models part I - A discussion
23 of principles *Journal of Hydrology*, vol. 10, no. 3, pp. 282-290, 1970.
- 24 [Nicodemus, K. K. Letter to the editor: On the stability and ranking of predictors from random](#)
25 [forest variable importance measures. *Briefings in Bioinformatics*, 12:369–373, 2011. doi:](#)
26 [10.1093/bib/bbr016](#)
- 27 [Nicodemus, K. K., Malley, J. D., Strobl, C., and Ziegler, A. The behaviour of random forest](#)
28 [permutation-based variable importance measures under predictor correlation. *BMC*](#)
29 [Bioinformatics](#), 11:110, 2010. doi: 10.1186/1471-2105-11-110

- 1 Palczewska, A., Palczewski, J., Robinson, R.M. and Neagu, D.: Interpreting random forest models
2 using a feature contribution method, Proceedings of the 2013 IEEE 14th International
3 Conference on Information Reuse and Integration, IEEE IRI 2013, pp. 112, 2013.
- 4 [Pappenberger, F., Matgen, P., Beven, K.J., Henry, J.B., Pfister, L., de Fraipont, P. Influence of
5 uncertain boundary conditions and model structure on flood inundation predictions. *Adv.
6 Water Resour.* 29, 1430–1449, 2006. Doi: 10.1016/j.advwatres.2005.11.012](#)
- 7 Pappenberger, F., Dutra, E., Wetterhall, F. & Cloke, H.L.: Deriving global flood hazard maps of
8 fluvial floods through a physical model cascade, *Hydrology and Earth System Sciences*, vol.
9 16, no. 11, pp. 4143-4156, 2012.
- 10 Peel, M. C., Finlayson, B. L., and McMahon, T. A.: Updated world map of the Köppen-Geiger
11 climate classification, *Hydrol. Earth Syst. Sci.*, 11, 1633-1644, doi: 10.5194/hess-11-1633-
12 2007, 2007.
- 13 Pelletier, P.M.: Uncertainties in the single determination of river discharge: a literature review.
14 *Canadian Journal of Civil Engineering*, 15:834–850, 1988.
- 15 [Rosso, R. A linear approach to the influence of discharge measurement error on flood estimates.
16 *Hydrol. Sci. J.* 30 \(1\), 137–149, 1998. doi: 10.1080/02626668509490975](#)
- 17 [Sandri, M. and Zuccolotto, P.](#) A bias correlation algorithm for the Gini variable importance
18 measure in classification trees. *Journal of Computational and Graphical Statistics*, 17:611-
19 628, 2008. [184], doi: 10.1198/106186008X344522
- 20 Schumann, Guy, Paul D. Bates, Matthew S. Horritt, Patrick Matgen, and Florian Pappenberger.:
21 Progress in Integration of Remote Sensing–derived Flood Extent and Stage Data and
22 Hydraulic Models. *Reviews of Geophysics* 47, RG4001, no. 4, 2009.
23 doi:10.1029/2008RG000274.
- 24 South African Water Affairs (DWA) database, <http://www.dwa.gov.za/Hydrology/>, last access: 10
25 July 2013.
- 26 Strobl, C., Malley, J. & Tutz, G.: An Introduction to Recursive Partitioning: Rationale,
27 Application, and Characteristics of Classification and Regression Trees, Bagging, and
28 Random Forests, *Psychological methods*, vol. 14, no. 4, pp. 323-348, 2009. doi:
29 10.1186/1471-2105-9-307

- 1 Thielen, J., Bartholmes, J., Ramos, M.-H., and de Roo, A.: The European Flood Alert System. Part
2 1: Concept and development, *Hydrol. Earth Syst. Sci.*, 13, 125–140, doi: 10.5194/hess-13-
3 125-2009, 9278, 2009.
- 4 Thiemig, V., Bisselink, B., Pappenberger, F., and Thielen, J.: A pan-African Flood Forecasting
5 System, *Hydrol. Earth Syst. Sci. Discuss.*, 11, 5559-5597, doi:10.5194/hessd-11-5559-2014,
6 2014.
- 7 [Tolosi, L. and Lengauer, T. Classification with correlated features: unreliability of feature ranking
8 and solutions. *Bioinformatics*, 27:1986–1994, 2011. doi: 10.1093/bioinformatics/btr300](#)
- 9 Tomkins, K.M.: Uncertainty in streamflow rating curves: Methods, controls and consequences.
10 *Hydrological Processes*, 28(3), pp. 464-481, 2014.
- 11 UNISDR: Global Assessment Report: Revealing Risk, Redefining Development, Chapter 2.2.
12 Global disaster risk trends, United Nations, printed in the UK, ISBN 978-92-1-132030-5,
13 page 22-27, 2011.
- 14 [UNOSAT, UNITAR Operational Satellite Applications Programme
15 <http://www.unitar.org/unosat/maps>, last accessed 1 December 2013](#)
- 16 Van Westen, C.J.: Remote sensing and GIS for natural hazards assessment and disaster risk
17 management. In: Shroder, J. (Editor in Chief), Bishop, M.P. (Ed.), *Treatise on
18 Geomorphology*. Academic Press, San Diego, CA, vol. 3, Remote Sensing and GIScience in
19 *Geomorphology*, pp. 259–298, 2013.
- 20 Yamazaki, D., O'Loughlin, F., Trigg, M.A., Miller, Z.F., Pavelsky, T.M. & Bates, P.D. 2014,
21 "Development of the global width database for large rivers", *Water Resour. Res.*, 50, 3467–
22 3480, doi: 10.1002/2013WR014664, 2014.
- 23 [Yitzhaki, S., Schechtman, E. *The Gini Methodology. A Primer on a Statistical Methodology*. 2013.
24 *Springer Series in Statistics. Volume 272*, 2013, ISBN: 978-1-4614-4720-7.](#)
- 25 Zaraj, Z., Zambrano-Bigiarini, M., Salamon, P. Burek, P., Gentile, A., and Bianchi, A.: Calibration
26 of the LISFLOOD hydrological model for Europe. Calibration Round 2013JRC Technical
27 Report, European Commission, Joint Research Centre. Ispra, Italy, 2013.
- 28 Zhang, Y., Hong, Y., Wang, X., Gourley, J.J., Gao, J., Vergara, H.J. and Yong, B.: Assimilation
29 of passive microwave streamflow signals for improving flood forecasting: A first study in
30 Cubango River Basin, Africa. *IEEE Journal of Selected Topics in Applied Earth
31 Observations and Remote Sensing*, 6(6), pp. 2375-2390, 2013.

1 Zhu, Z., Bi, J., Pan, Y., Ganguly, S., Anav, A., Xu, L., Samanta, A., Piao, S., Nemani, R.R. &
2 Myneni, R.B.: Global data sets of vegetation leaf area index (LAI)3g and fraction of
3 photosynthetically active radiation (FPAR)3g derived from global inventory modeling and
4 mapping studies (GIMMS) normalized difference vegetation index (NDVI3G) for the period
5 1981 to 2011, Remote Sensing, vol. 5, no. 2, pp. 927-948, 2013.

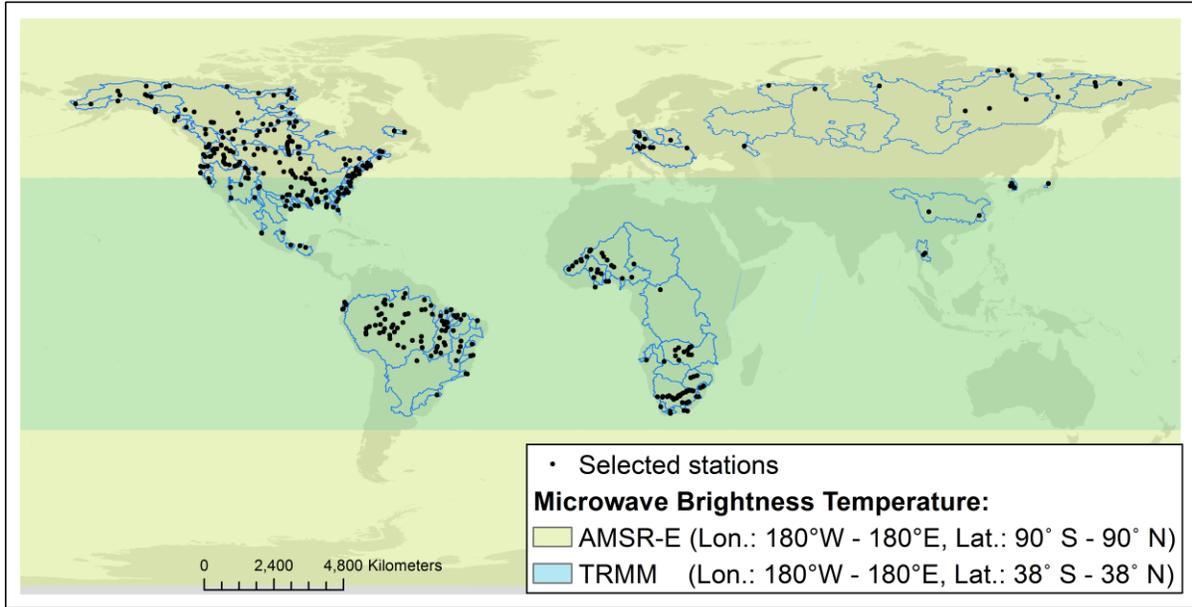
1 Table 1 Number of catchments by continent and range of upstream areas for the located
 2 stations.¹Stations used for calibration and validation.² South African upstream catchment areas are
 3 not available.

Continent	Number of satellite locations for extraction (n=398)	Number of stations for calibration(n=322)	Number of Catchment¹	Upstream catchment (km²) Approx. range
Africa	75	51	21	46990 – 850500 ²
Asia	23	3	4	7150 - 11000
Europe	13	7	3	9000 - 132000
North America	207	183	86	5300 - 1850000
South America	80	78	38	1400 - 4680000

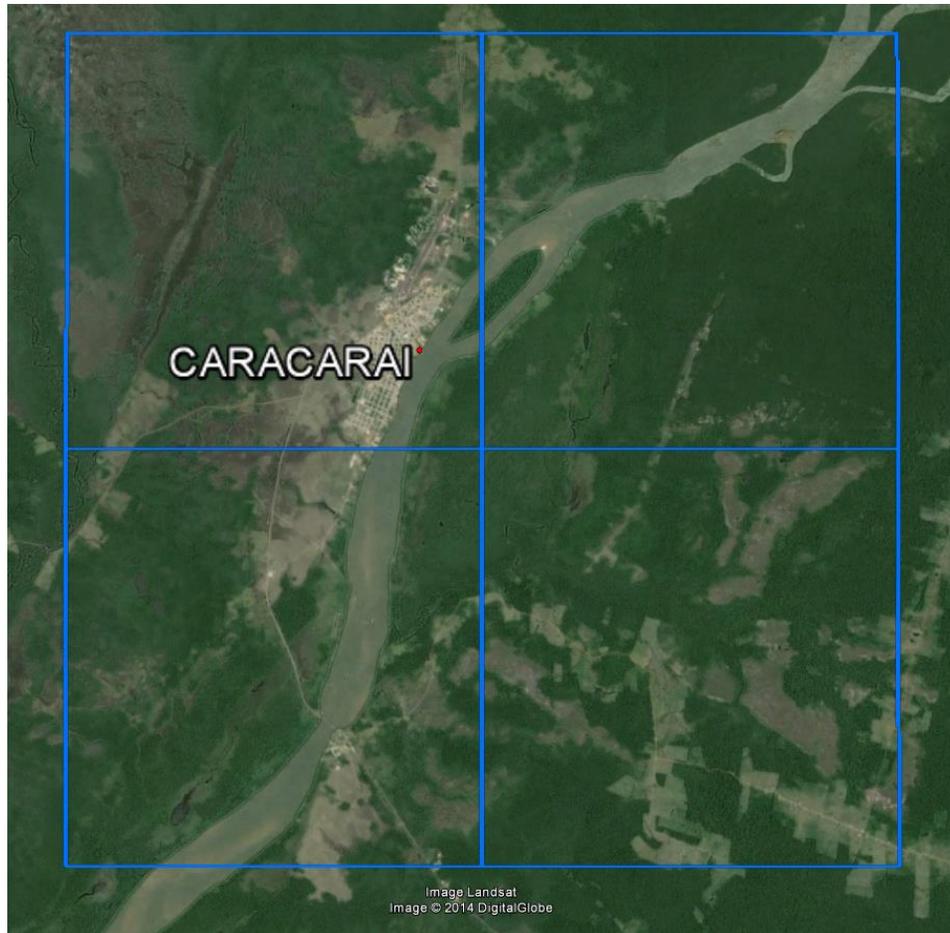
4

1 Table 2 Climate and land cover type of the 322 sites selected for the calibration and validation,
 2 aggregated by continent, climate, and land cover. ¹ Vegetation means a combination of grassland,
 3 shrubland and forest. ²Types of land cover and climate where the number of locations in each type
 4 was very low (e.g. 3) were excluded for their respective variables analysis as they will not be
 5 representative on a global scale.

Climate	Africa	Asia	Europe	North America	South America	Total
Arid	30			25		55
Tropical	10				75	85
Temperate	11		3	51	3	68
Cold		3	4	104		111
Polar ²				3		3
Total	51	3	7	183	78	322
Land cover	Africa	Asia	Europe	North America	South America	Total
Open Forest	4			23		27
Closed to Open Forest	16	1	1	16	41	75
Closed Forest				33		33
Mosaic Vegetation predominant ¹	19	2		47	24	92
Mosaic cropland or grassland predominant	5		1	26	9	41
Rainfed crop			4	5	4	13
Sparse vegetation	2			14		16
Sparse vegetation+crops	5			8		13
Urban			1	10		11
Bare areas ²				1		1
Total	51	3	7	183	78	322

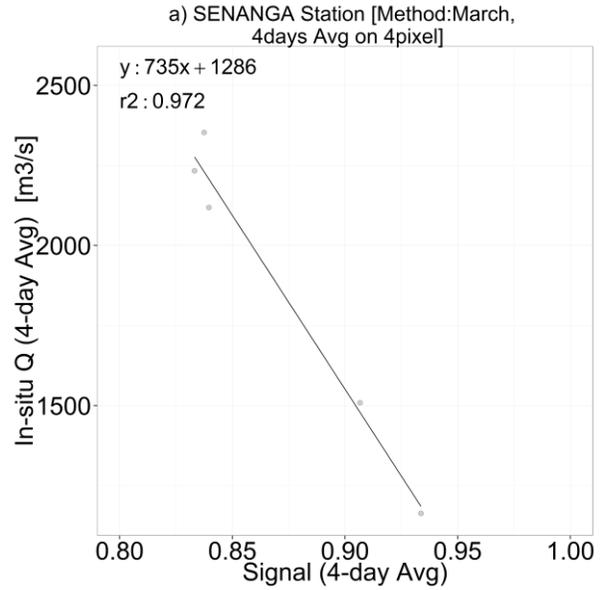


1
 2 Figure 1 Location of selected stations (398) and corresponding river basins (109). TRMM and
 3 AMSR-E brightness temperature product extents are also provided.

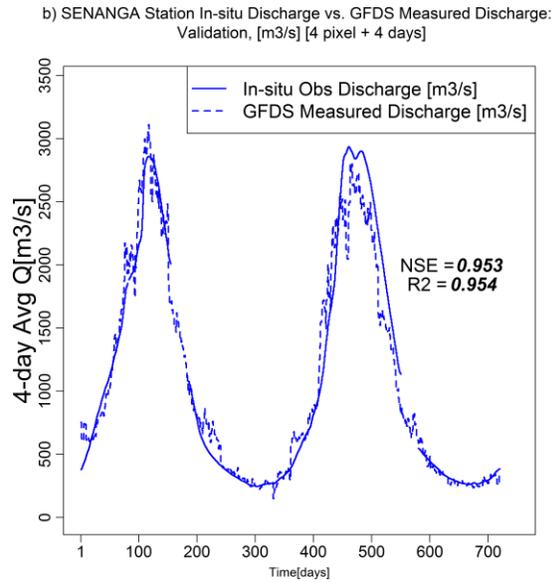
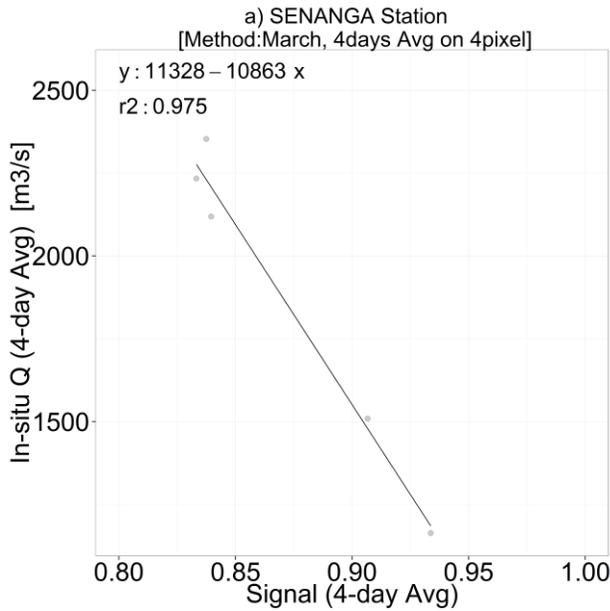


1
2 Figure 2 Example of a measurement site: Caracarai station (Rio Branco Catchment, Brazil). The
3 blue rectangles outline the measurement pixels and background image is from 2014 Google
4 (Landsat, DigitalGlobe).

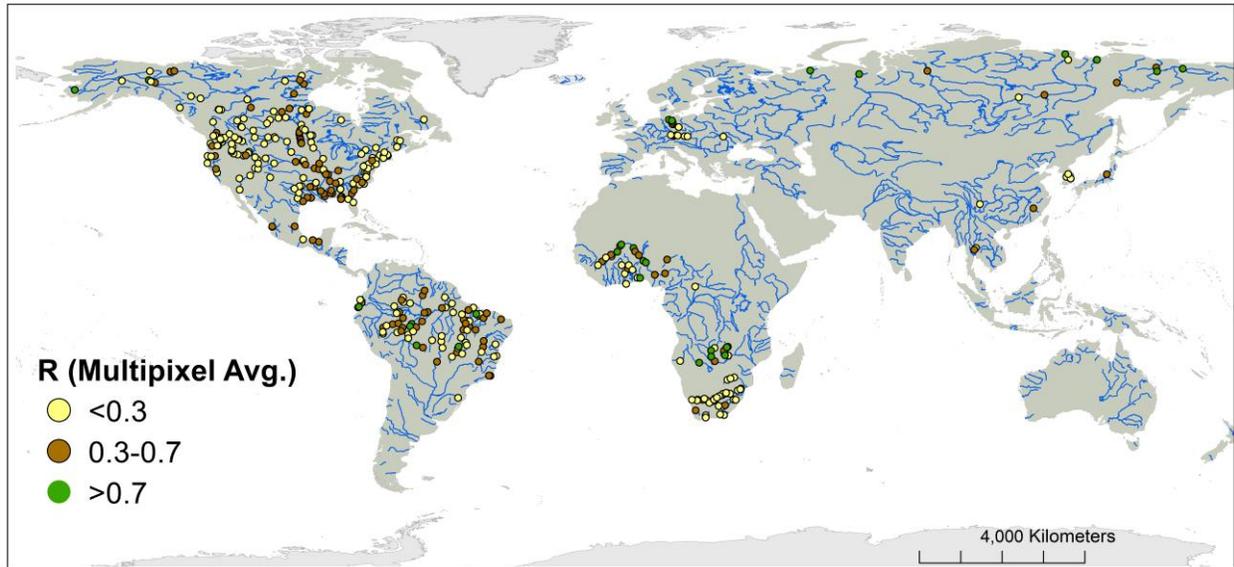
1



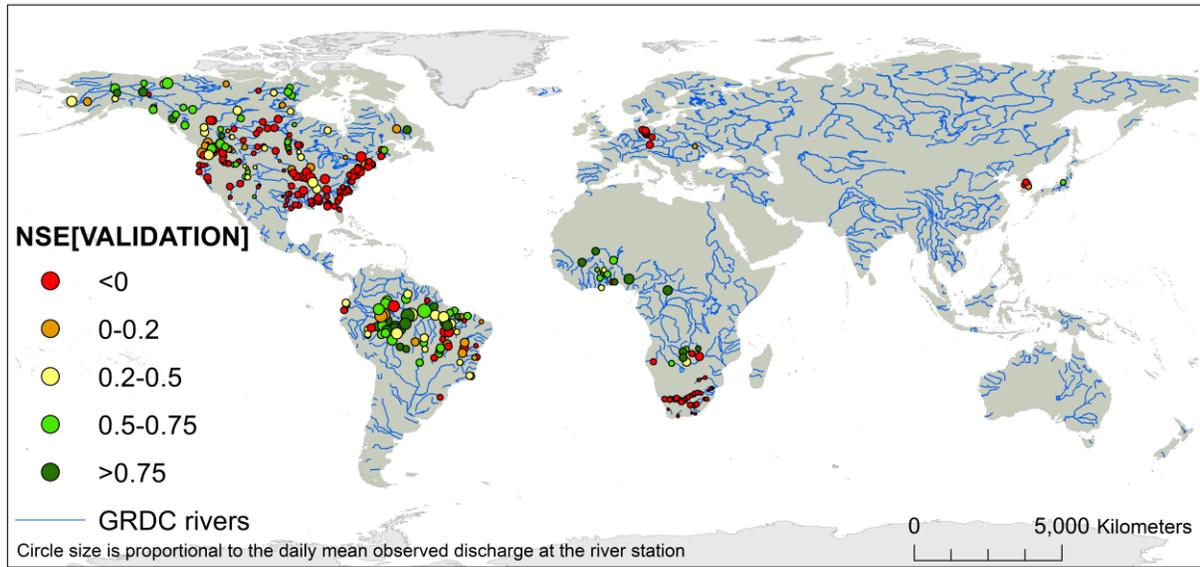
2



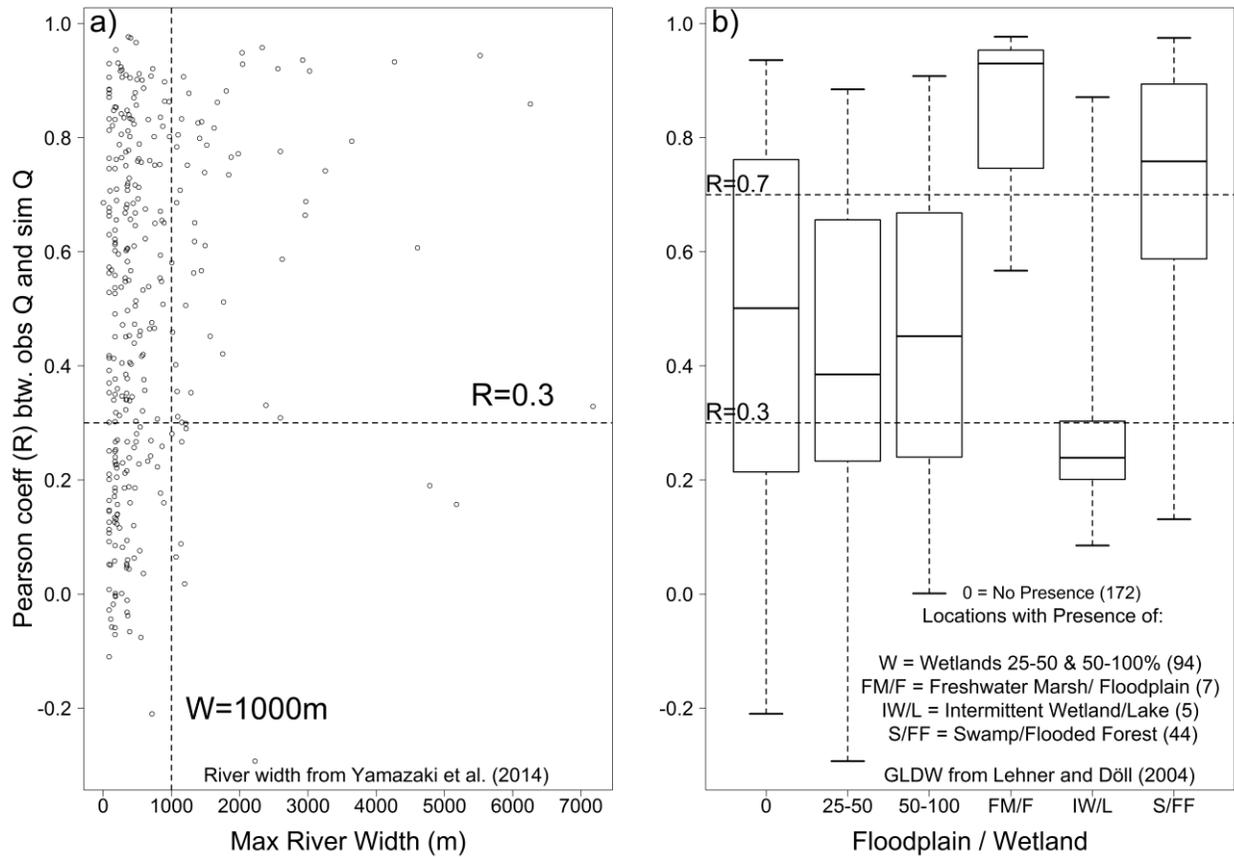
3 Figure 3 a) scatterplot for the Senanga station (Long 23.25, Lat. -16.116) in the Zambezi River
4 (Africa). Monthly mean for March from 1998 up to 2002. b) Validation hydrograph for 2003-2004
5 and skill scores for Senanga.-The (monthly) rating equations were used to calibrate the signal into
6 discharge units. [Different rating equations were used for different months.](#)



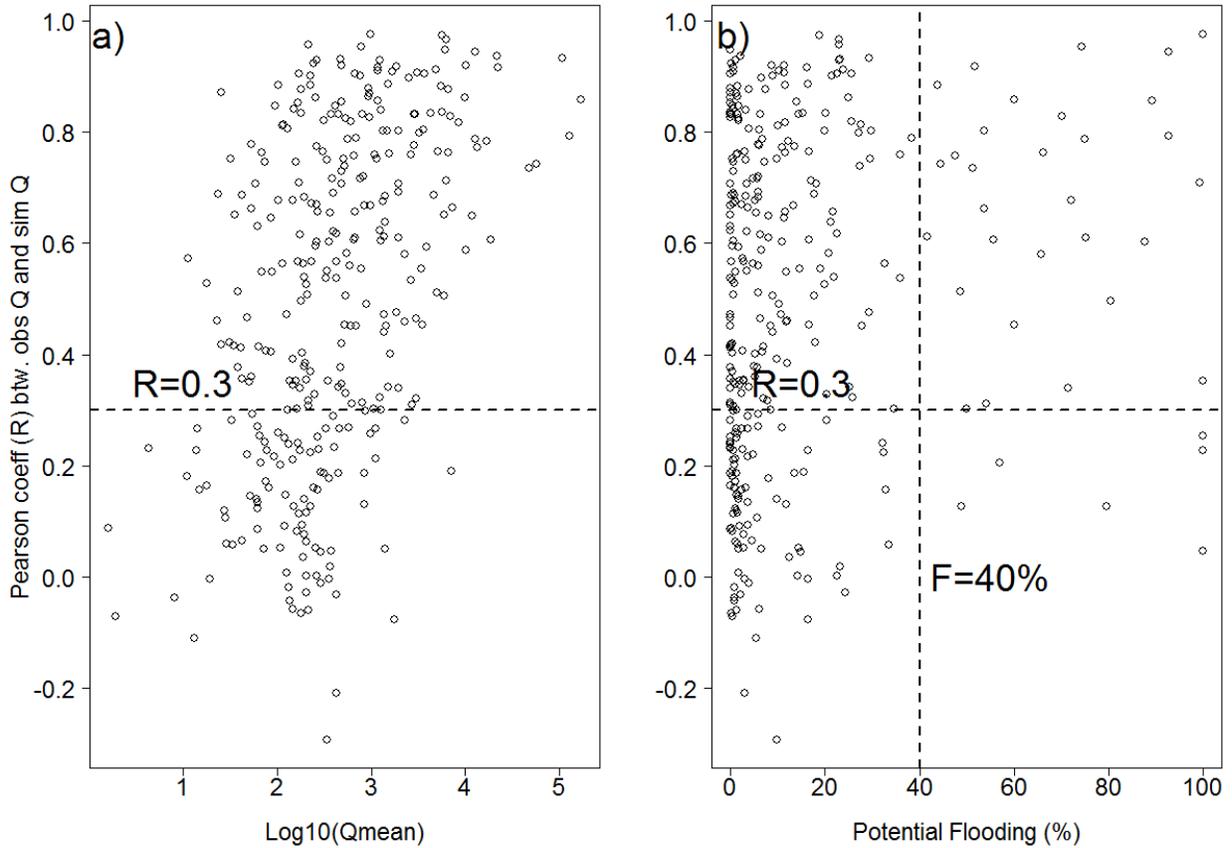
1
2 Figure 4 Location of stations and R skill score between in situ observed discharge and satellite
3 signal (4 days and 4 pixels average). Globally, 169 stations have $R > 0.3$.



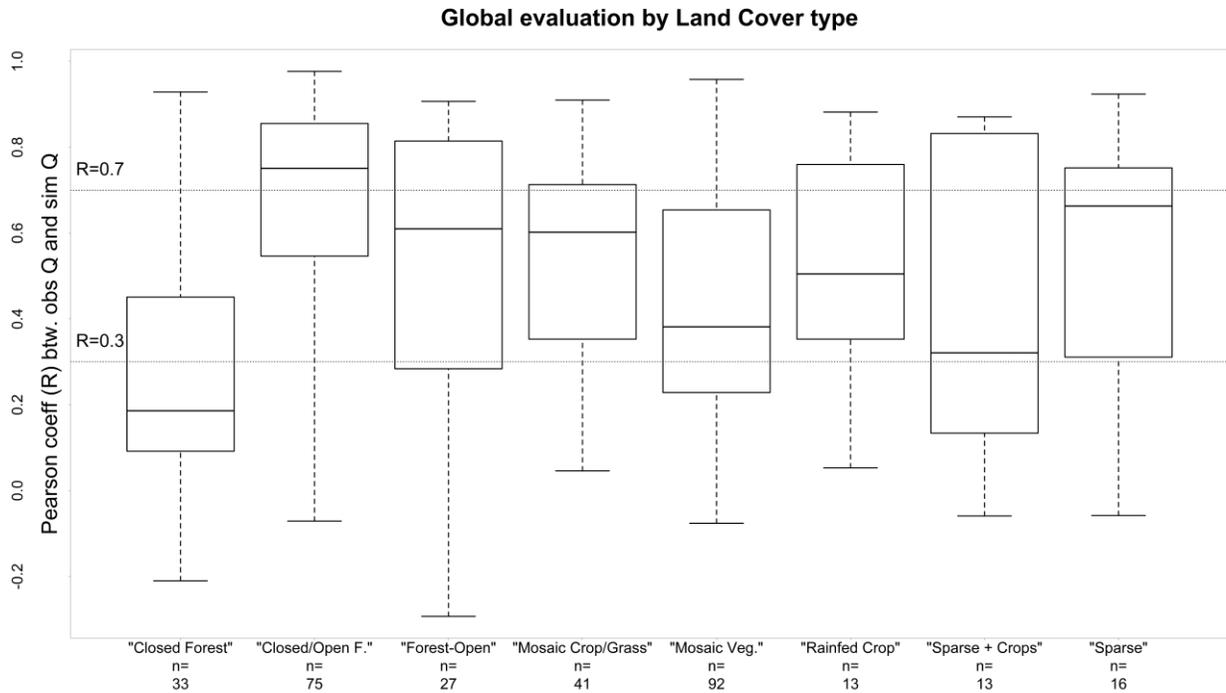
- 1
- 2 Figure 5 Nash-Sutcliffe efficiency of the validation (n= 332 stations). Globally, 154 stations have
- 3 NSE>0 of which 80 stations have NSE> 0.50.



1
 2 Figure 6 a) relationship between R obtained from the validation of satellite measured discharge and
 3 the maximum river width for each location; b) relationship between the same R score and the
 4 presence of significant floodplains, flooded forest and wetlands Horizontal dotted line shows the
 5 $R=0.3$ and $R=0.7$ threshold, the vertical line is the river width equal to 1km.

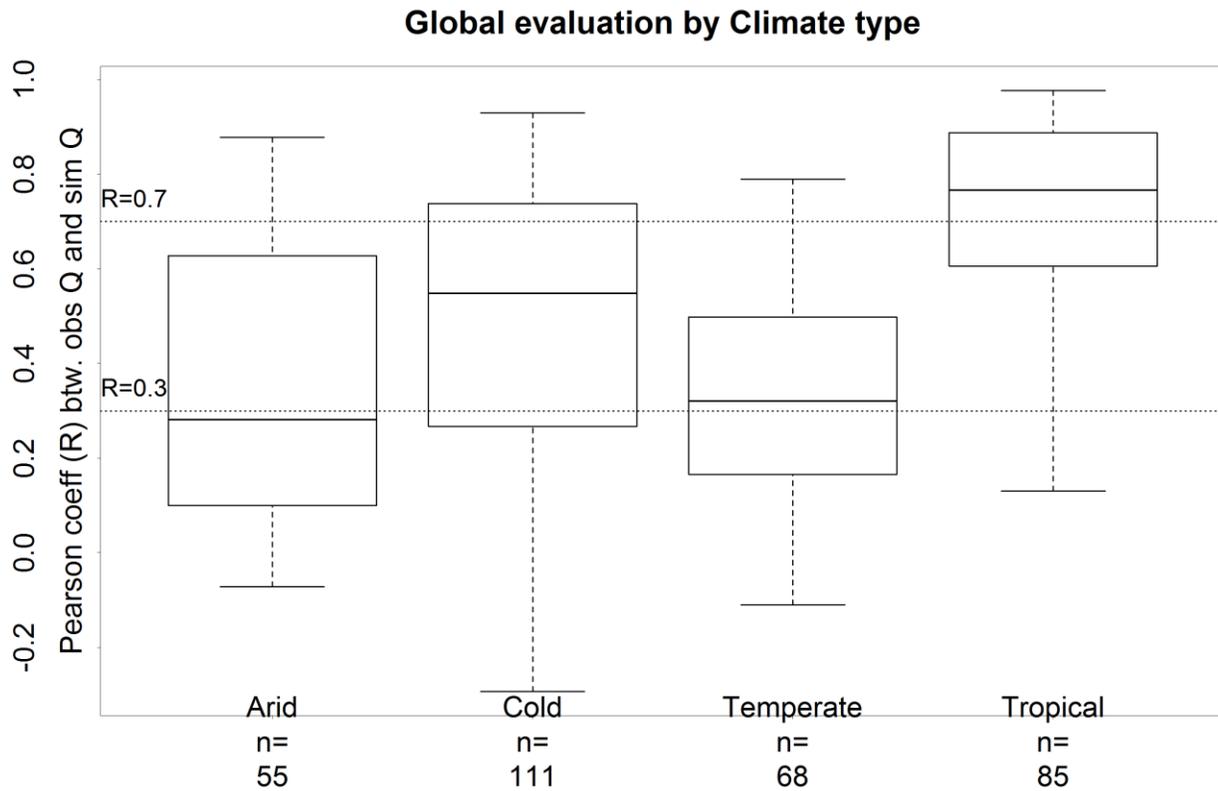


1
 2 Figure 7a) relationship between R obtained from the validation of satellite measured discharge and
 3 the mean in situ observed discharge (log10 displayed) for each station; b) relationship between the
 4 same R score and the potential percentage of flooded area per pixel for a 100 year return period
 5 flood event (Pappenberger et al., 2012). Horizontal dotted line shows the R=0.3 threshold, the
 6 vertical line is the 40% potential flooding threshold.



1

2 Figure 8 Global evaluation of the R score obtained during the validation and its classification by
 3 the land cover type of the stations. Land cover type were aggregated from the GlobCover (2009)
 4 and modified by visual check with Google maps. Note that artificial and bare land cover were
 5 excluded on this figure.

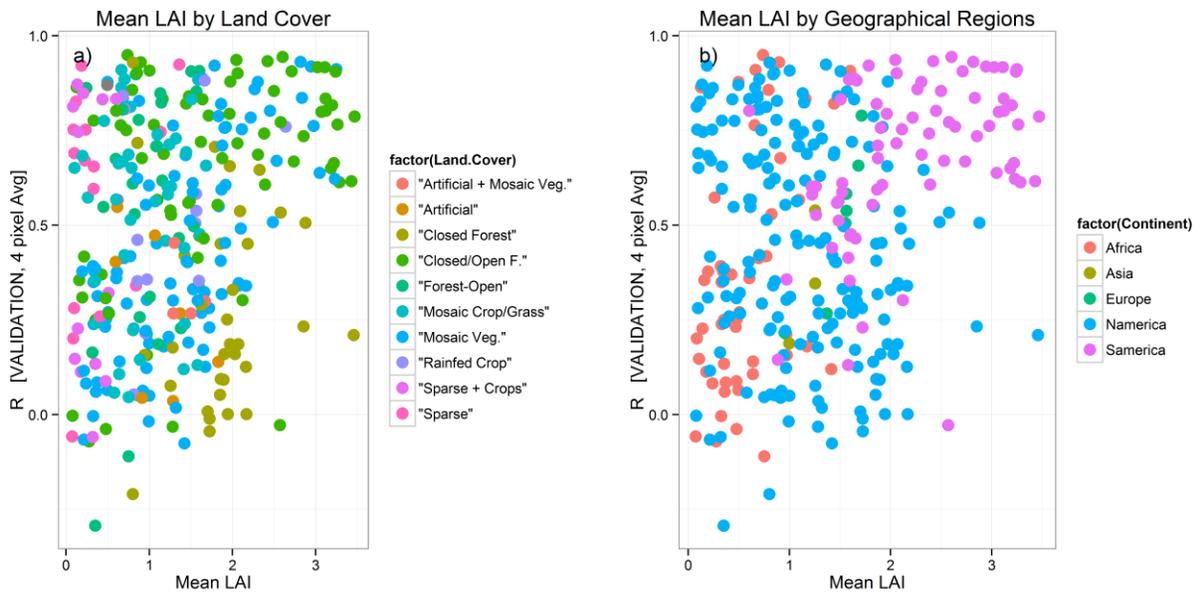


1

2 Figure 9 Global evaluation of the R score obtained during the validation and its classification –only

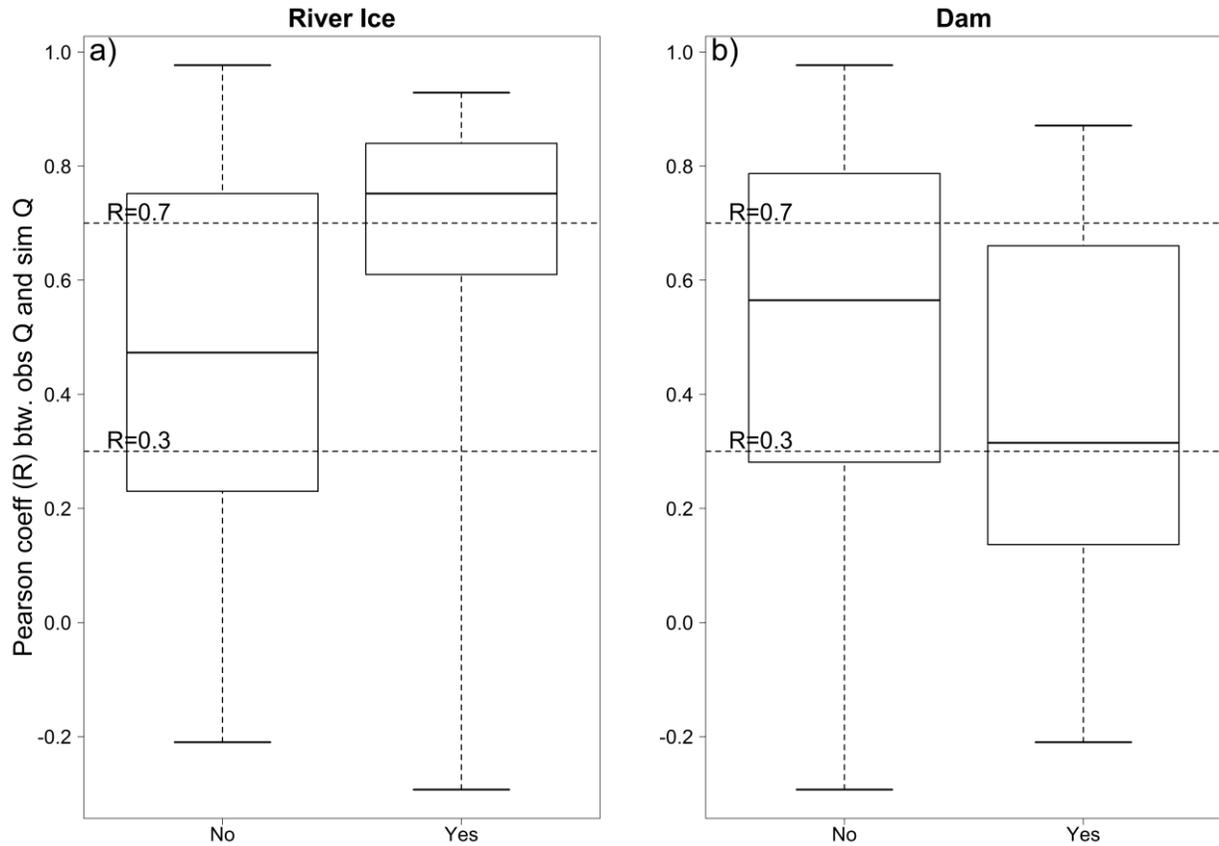
3 main types-by the Köppen-Geiger climate area (Peel et al. 2007). Note that polar climate was

4 excluded from this analysis as only three stations fell into this category.

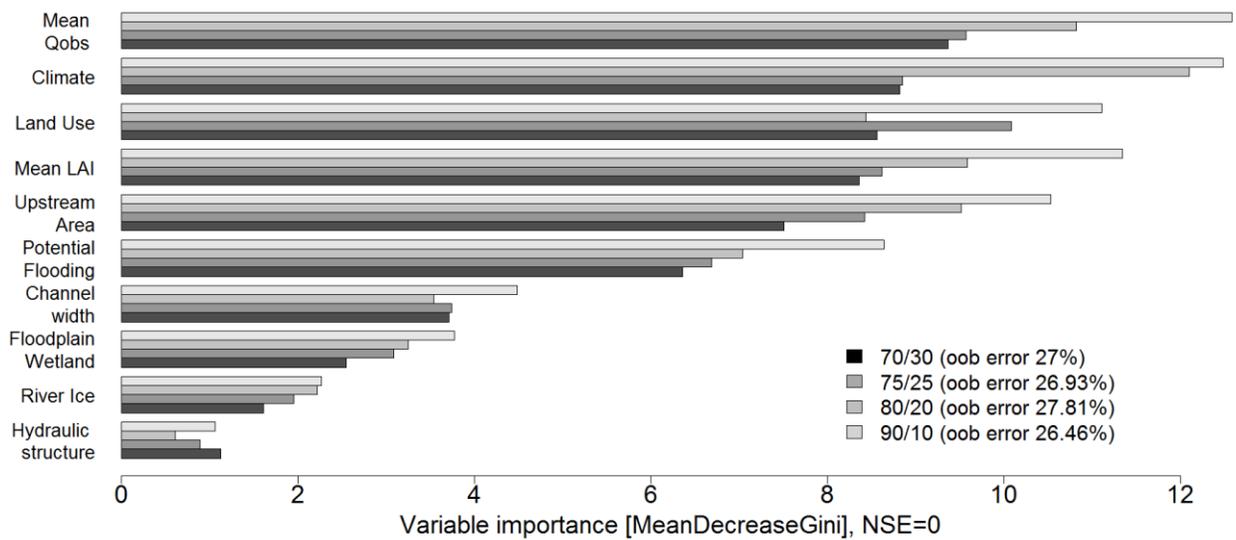


1

2 Figure 10 Evaluation of the R score obtained during the validation and its classification by Leaf
 3 Area Index (LAI), also a factor of land cover and geographical regions.



1
 2 Figure 11 Evaluation of the R score obtained during the validation and its classification by a)
 3 presence or not of a river ice (Brown et al., 2002), b) presence or absence of a nearby dam or
 4 hydraulic control infrastructure using the Global Reservoir and Dam (GRanD) (Lehner et al.,
 5 2008) and visual check from Google maps.. To note that for the validated locations, all stations with
 6 river ice and most of them with dams and are located in North America.



1

2 Figure 12 Average variable importance of 200 runs using the Random Forest methodology. Nash-
 3 Sutcliffe score was chosen as a quality index to categorised the stations as true (good predictive)
 4 or false the stations as false (poor predictive). With a threshold of NSE=0, we have about 50% of the
 5 stations above and below that value. Results are shown for the different training and test groups.
 6 For all the test groups and runs, the average highest variable importance was obtained for mean
 7 observed discharge, climatic region, land cover/ mean LAI and upstream catchment area, and the
 8 lowest for dam/hydraulic structure presence and river ice.

1 **Appendix A: Land cover types**

2 Table A 1 Studied land cover types from GlobCover (2009) aggregated into broader categorical
 3 classes by type and vegetation density.

Label	Aggregated classes
Rainfed croplands	Rainfed croplands
Sparse (<15%) vegetation	Sparse vegetation
Closed to open (>15%) broadleaved evergreen or semi-deciduous forest (>5m)	
Closed to open (>15%) mixed broadleaved and needleleaved forest (>5m)	
Closed to open (>15%) (broadleaved or needleleaved, evergreen or deciduous) shrubland (<5m)	
Closed to open (>15%) herbaceous vegetation (grassland, savannas or lichens/mosses)	Closed to open forest
Closed to open (>15%) broadleaved forest regularly flooded (semi-permanently or temporarily) - Fresh or brackish water	
Closed to open (>15%) grassland or woody vegetation on regularly flooded or waterlogged soil - Fresh, brackish or saline water	
Open (15-40%) broadleaved deciduous forest/woodland (>5m)	Open forest
Open (15-40%) needleleaved deciduous or evergreen forest (>5m)	
Mosaic cropland (50-70%) / vegetation (grassland/shrubland/forest) (20-50%)	Mosaic cropland or grassland
Mosaic grassland (50-70%) / forest or shrubland (20-50%)	
Mosaic vegetation (grassland/shrubland/forest) (50-70%) / cropland (20-50%)	Mosaic vegetation predominant
Mosaic forest or shrubland (50-70%) / grassland (20-50%)	
Closed (>40%) broadleaved deciduous forest (>5m)	
Closed (>40%) needleleaved evergreen forest (>5m)	
Closed (>40%) broadleaved forest or shrubland permanently flooded - Saline or brackish water	Closed forest
Artificial surfaces and associated areas (Urban areas >50%)	Urban

