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

**B. Revilla-Romero [1,2], J. Thielen[1], P. Salamon[1], T. De Groeve[1], G. R. Brakenridge[3]**

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

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

[3] University of Colorado, Boulder, USA

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

## Abstract

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

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

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

33

## 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 Pilotand 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 Mission[1] (TRMM) operational from 1998 to present and the Advanced

81  Microwave Scanning Radiometer for Earth Observation System[2] (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-W1[3], 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/flooddetection/) 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

## 2   Study regions and data

125

### 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), Le

144     Coz et al. (2014), and Tominsk (2014). However, the uncertainty of river discharge is even higher

145　during floods events when the stage-discharge relationship, the so-called rating curve, is used. As

146　evaluated by Pappenberger et al. (2006), the analysis of rating curve uncertainties leads to an

147　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 Mackezie 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

## 3   Methodology

235

### 3.1. Satellite signal extraction

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

(INSERT FIGURE 2 HERE)

255

### 3.2. Satellite signal calibration and validation

For those co-located ground stations and satellite measurement sites where both sets of data (signal 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.

$$Q_{\text{GFDSmeasured of X month}} = a_{\text{month}} + b_{\text{month}} * \text{signal} \tag{1}$$

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. Power law fitting was also tested to calibrate the signal into discharge units yielding

287   similar results (see Open Discussion Author's Response).

288 The validation of the satellite derived daily discharge data was carried out with daily in situ data
289 on a two-year period, and skills scores were calculated to quantify the agreement between both
290 satellite and ground measured discharge. We are aware of the limited number of years (data) with
291 available time series for both variables, which might influence the robustness of the calibration.
292 In some cases there were longer time series available, but to standardised the analysis for all the
293 stations we used five years (1998-2002 or 2003-2008 for Northern stations with AMSR-E signal)
294 and the following two years for validation purposes (2003-2004 and 2009-2010 respectively). Note
295 that for 36 out of the 322 stations available data length was between six years and three months to
296 almost seven years. Validation was still carried out for the same period, but the data used for
297 calibration was slightly reduced. As an example, Figure 3a presents the scatterplot for the month
298 of March for the Senanga Station (Long 23.25, Lat. -16.116) in the Zambezi River (Africa) with
299 mean values derived from the period 1998 to 2002. For the same location, Figure 3b shows the in
300 situ observed and the GFDS measured discharge derived from the GFDS signal for the period
301 2003-2004.

302 (INSERT FIGURE 3 HERE)

303

304 **3.3. Skill scores**

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

310 One of the advantages of the R coefficient is its independence on the units of measurement, which
311 permits the comparison of dimensionless GFDS signal data. A small value indicates a weak or
312 non-linear relationship between the satellite signal and discharge. For this study, we grouped the
313 computed R values into three ranges as follows: <0.3, [0.3-0.7], and >0.7. While Pearson
314 benchmarks linear relationship, Spearman benchmarks monotonic relationship. Spearman's
315 validation scores just obtained a mean value 6% higher than Pearson mean score (see Open
316 Discussion Author's Response). On this manuscript, results are analysed based on the scores
317 obtained using Pearson correlation coefficient.

318  Nash-Sutcliffe efficiency (NSE) (Nash and Sutcliffe, 1970) is typically used to assess the
319  predictive power of hydrological models and was here calculated to describe the accuracy of
320  satellite-derived discharge in comparison to gauge-observed discharge values. Higher values of
321  the Nash-Sutcliffe statistic should indicate more correlated results, without other factors taken into
322  account, such as autocorrelation (Brakenridge et al., 2012). However, the degree of correlation of
323  these variables does not verify the discharge magnitudes (Brakenridge et al., 2013). A NSE value
324  of 1 corresponds to a perfect match of modelled to the observed data whereas NSE = 0 indicates
325  that the model predictions are as accurate as the mean of the observed data. The resulting scores
326  will be classified as in Zaraj, et al. (2013): $< 0$, [0.2-0.5], [0.5-0.75], and $> 0.75$.

327

328  **3.4 Factors affecting the satellite signal**

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

335  First, we use the skill scores (R and NSE) obtained from a simple analysis for each individual
336  factor or variable. Second, we seek to understand which of the surface variables have the greatest
337  importance in determining sites with a good or poor performance. For this purpose, we use a
338  decision tree technique called Random Forest (RF). Among other features, this allows ranking of
339  the relative importance of each variable. The technique is described by Breiman (2001) and
340  implemented in R by Liaw and Wiener (2002), where the reader is referred for a more detailed
341  explanation. As a summary of the Random Forest algorithm, *ntree* bootstrap samples are randomly
342  selected from the data set, a different subset is used for each bootstrap and for each sample a tree
343  is grown, obtaining *ntree* trees. Random Forest is called an ensemble method because it applies
344  the method for a number of decision trees, in this case 500, in order to improve the classification
345  rate. Some stations are left out of the sample (out-of-bag) and used to gain an internal unbiased
346  estimate of the generalisation error (oob errors) and to obtain estimates of the importance of the
347  variables (Breiman, 2001). These values are averaged over the *ntree* trees. For the variables

classification, the node impurity is measured by the Gini index. 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 (Gonzalez,L., et al. 2010).. The variables with higher decrease in Gini values (lower Gini) are those with higher importance on the classification analysis.

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

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

## 4    Results and discussion

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

378     discharge, river width, upstream catchment area, potential flood hazard area, land cover, leaf area

379     index, climatic zones, presence of large floodplains, flooded forest and wetlands, river ice and

380     hydrologic structure. Finally, the relative importance of all variables in comparison to each other

381     has been assessed (Section 4.4).

382     Before analysing the validation results, it is important to highlight two possible different sources

383     of error which might influence the outputs. Firstly, the signal to noise ratio might be low for a site

384     or have intermittent instrument noise occasionally producing positive spikes in discharge.

385     Secondly, the rating curve may be offset, which will result in a consistent bias on the discharge

386     values for that location even though the time series are strongly correlated.

387

## 4.1. Correlation of raw satellite data vs. gauge observations

389     The first step was to look at the "raw" correlation between daily ground station-measured water

390     discharge and the satellite signal and to calculate the empirical linear relation between these two

391     variables for each site. The full time series, including low flows, were used for the calculation and

392     executed for 398 stations. Figure 4 shows the R skills obtained. 169 out of 398 sites have an R >

393     0.3 and 42 of them have R>0.5. Perhaps, correlations might have been higher if regression would

394     have not been restricted to linear equations (Brakenridge et al., 2007, 2012).

395     (INSERT FIGURE 4 HERE)

396

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

398     For the stations with over six years of contemporary data for both in situ discharge and satellite

399     signal, we obtained regression equations for each month of the year and station using the first five

400     years of data. Next, using these equations we carry out a calibration of the daily signal into

401     discharge units.  Afterwards, the validation of the GFDS measured discharge was implemented for

402     the following two years. In some regions such as Northern Asia, the lack of available recent long

403     time series (after 2002) meant that the number of stations available for calibrating the satellite into

404     discharge measurements was reduced.  Stations where the number of years matching observed

405     discharge and satellite signal was shorter than six years were excluded from the validation exercise

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

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

414 (INSERT FIGURE 5 HERE)
415

416 **4.3. Analysis of the factors affecting the satellite signal**
417

418 ### 4.3.1. River width and presence of floodplain and wetlands.
419 As a first step to analyse the potential relationship between the individual local characteristics and
420 the performance of the locations in global terms, we study the R score of the validation for the 322
421 stations in relation with the maximum river width value at each location (Figure 6a). Results
422 indicate that locations with a river width higher than 1 km are more likely to score an R larger than
423 0.3. In fact, the mean R score is 0.60. Where 26 out of 64 (~41%) have R> 0.75. However, there
424 is a number of stations with lower river width that also obtained high scores. As the retrieval of
425 the satellite signal also depends on the floodplain geometry. As soon as the river floods and water
426 goes over-bank, the proportion of water in the wet pixel greatly increases. So the score should be
427 also high for small rivers with a proportionally big floodplain. Figure 6b shows the R scores by
428 locations where the majority of the area belongs to floodplain, flooded forest and wetlands
429 category or, their absence. In our study, higher median scores were obtained for those located in
430 large freshwater marsh and floodplains, followed by those on swamps and flooded forest. These
431 results give a first indication on the characteristics of the locations with better performance.

432 (INSERT FIGURE 6 HERE)

433 ### 4.3.2. River discharge and potential flooding
434 Flooding is determined by the discharge as well as the potential flood hazard. Figure 7a shows that
435 84 out of 95 stations with R<0.3, also have mean discharge values lower than 500 $m^3s^{-1}$ (Log10

436  (500) ≈2.7), of which 55 stations in fact had a mean discharge lower than 200 $m^3s^{-1}$. These stations

437  are mainly located in South Africa, and in some areas of North America. Therefore, it can be

438  concluded that the mean discharge can be considered a key variable that determines the

439  appropriateness of locations for which satellite discharges can be derived: As 77% of the stations

440  with Q<500 $m^3$/s, have R< 0.3, while 91.5% of the stations with Q>500 $m^3$/s have R >0.3, locations

441  with discharge of less than 500 $m^3s^{-1}$ might not provide reliable results for a global satellite-based

442  monitoring system. Alternatively, non-permanent rivers and streams exhibiting only seasonal or

443  ephemeral flow (typical for dry regions) may require a different monitoring approach, wherein a

444  "dry" threshold is established for the signal data.

445  After excluding the global stations with low skill score due to low flows and studying the

446  remaining stations, we can better understand the performance of the system in relation to other

447  local characteristics. Figure 7b shows for each location the relationship between the validation R

448  and the percentage of area in each pixel covered by potential flooding during a 100 year return

449  period flood event, obtained with the model grid (HTESSEL+CaMa-Flood) (downscaled from a

450  25 x 25 km pixel, Pappenberger et al., 2012). 100 means totally flooded across its area, 50 means

451  50 % of the area within the cells is flooded, and 0 means that the area is not flooded. Although

452  there is not a clear trend for all the points, result indicate that locations with a percentage of

453  potential flooding larger than 40%, are expected to score an R larger than 0.3.

454  (INSERT FIGURE 7 HERE)

455  ### 4.3.3. Land cover types and climatic areas

456  Figure 8 presents a global evaluation of the R score obtained during the validation and its

457  classification by the land cover type of the stations. The bare land cover category was excluded

458  from this study as only one of the selected locations belong to that class. Looking at the median of

459  the boxplot (see Figure 8), we found that some of the locations with higher density of vegetation

460  such as those located on "closed forest" and "mosaic with predominant vegetation" (included

461  forest, scrublands and grasslands) obtained lower median scores values. In contrast, the locations

462  with lower vegetation density such as "sparse vegetation", "mosaics with predominant

463  cropland/grasslands", "open forest" and "closed to open forest" land cover types obtained larger

464  median R scores, around 0.6-0.8. Similar results can be observed when looking at the interquartile

465  range or spread of the boxplots: "closed to open forest" and "mosaics with predominant

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

473  (INSERT FIGURE 8 HERE)

474

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

483  (INSERT FIGURE 9 HERE)

484  ### 4.3.4.  Leaf Area Index (LAI)

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

494     We decided to study the relationship between the mean Leaf Area Index and the skill obtained in

495     the validation for each location, also looking at complementary variables such us the land cover

496     and the geographical region which the stations belong to. Figure 10 shows that locations with a

497     mean [LAI > 2] predominantly have a "closed to open forest" type in South America (31 stations)

498     of which 29 have an R score higher than 0.6. For [LAI > 2] there is also 12 North American

499     locations with "closed forest" land cover but in general with poorer scores for those locations.

500     Additionally, 18 stations with mosaic vegetation from North and South America obtained [LAI >

501     2] and 16 out of them, a [R>0.6]. For [LAI < 2], both the land cover and geographical locations

502     are distributed along the scatterplots, from poor to high correlations.

503     (INSERT FIGURE 10 HERE)

504

505     ### 4.3.5. River ice

506     Figure 11a shows the scores obtained for the locations with presence or absence of river ice,

507     including a range from continuous to sporadic (Brown et al., 2002). It can be seen that stations

508     located in areas with river ice tend to have a good correlation between in situ and  satellite

509     measured discharge (based on 33 stations), as the system tends to capture well the annual spring

510     ice break-up and freezing as indicated in the study by Brakenridge et al.(2007) and Kugler (2012).

511     At these locations, once ice-covered there is no sensing capability from the system: which may

512     seems analogous to low flow conditions, and for which sites we obtained lower scores. However,

513     there is an important difference when analysing time series of signal between ice covered high

514     latitude river and all-year-around low flow rivers. When on the sites with river ice melting process

515     takes place, there is an increase of runoff happening and for many places the signal strongly

516     indicates this increased flow. On the other type of rivers, low flows is generally a characteristic for

517     most of the year and if the signal to noise is low, the signal retrieved is very noisy: one motivation

518     for setting a "dry" threshold for such sites.

519

520     ### 4.3.6. Hydraulic structures

521     The correlation between satellite and discharge data depends on both variables. Typically it is

522     assumed that observed discharges are "ground truth", however, when influenced by structures and

523     dams the ground discharge may not be well-monitored by flow area/flow width variation. For

524     example, when there is a major increase in river discharge but a flood is avoided by artificial

525    levees, we cannot expect that the satellite signal will accurately capture the flood hydrograph; as

526    well, downstream flooding may be attenuated by an upstream flood control dam and reservoir; so

527    that the gauge location is critical. Figure 11b shows the influence of the presence or absence of a

528    nearby dam using the Global Reservoir and Dam (GRanD) database (Lehner et al., 2008) or

529    visually identified hydraulic control infrastructure. Locations where the dam or other element was

530    present (48 stations) obtained lower median R score. Therefore, ideally, observation sites should

531    be located in areas without hydraulic control infrastructures.

532    (INSERT FIGURE 11 HERE)

533

## 4.4. Variable importance

535    Based on the individual analysis of the signal potential influence factors we found that to

536    understand the site performances, on some occasions multiple variables need to be analysed in a

537    simultaneous way. For example, the general low scores obtained on the Eastern USA stations

538    might be due to a number of factors: ~64% of these stations have a mean discharge value lower

539    than 500 $m^3s^{-1}$ and ~88% of the stations are located at river width lower than 1km. In addition,

540    ~59% of the stations are located in wetlands areas.  Another example, in this case regarding the

541    exceptions of the low R and mean observed discharge higher than 500 $m^3s^{-1}$, all the 11 locations

542    have a potential probability of flooding lower than 21%, the land cover of 10 out of 11 is forest, 5

543    of them located in wetlands and two of them have a nearby hydraulic structure. Despite exhibiting

544    a mean discharge greater than 500 $m^3s^{-1}$, these other local characteristics may be the cause of the

545    poor performance.  Therefore, we decided to use a classification decision tree technique (Random

546    Forest), which split the dataset at each node according to the value of one variable at a time (the

547    best split) from a selected set of variables  to understand the importance of each variable. Random

548    Forest is called an ensemble method because it is performed for a number of decision trees, in this

549    case 500 trees, in order to improve the classification rate.

550    The result presented here is the rank of the importance of variables to classify a location with a

551    good or poor performance. These values are obtained as an output of the Random Forest analysis

552    and are, in addition, the average of 200 independent runs. As explained in section 3.4 the variable

553    importance based on the mean decrease in Gini index was calculated for the Nash-Sutcliffe

554    Efficiency (NSE) score obtained from the validation. We used a NSE=0 to distinguish the sites

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

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

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

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

578 (INSERT FIGURE 12 HERE)

579

580 **5   Conclusions and future research**

581

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

584    continental R skills are as follow: Africa 0.382, Asia 0.358, Europe 0.508, North America 0.451

585    and South America 0.694. Approximately 48% of these stations have an NSE score higher than

586    zero; 13 located in Africa, 77 in North America, 62 in South America, 1 in Asia and 1 in Europe.

587    Results showed that the majority of the stations that received low skills scores, were due to low

588    flow conditions. For example, 84 out of 95 stations with R<0.3, have mean discharge values lower

589    than 500 $m^3s^{-1}$. These are located mainly in South Africa with 25 cases and North America with

590    53 cases, which penalised their average continental skills. Note that our focus was on factors

591    affecting the method, globally, and that these skill values do not directly indicate at-a-site

592    measurement accuracy (which could be improved, for example, by use of non-linear rating

593    equations and/or accommodation of any phase shift or timing differences in flow area- versus state-

594    based discharge monitoring).

595    In order to better understand the impact of the local conditions on the performance of the sites, we

596    looked first at specific factors individually. In general terms, higher skill scores were obtained for

597    location with one or more than one of the following characteristics: a river width higher than 1km,

598    a large floodplain area, in flooded forest, with a potential flooded area per pixel greater than 40%

599    during a 100 year return period flood event, a land cover type of sparse vegetation, croplands or

600    grasslands and closed to open and open forest, Leaf Area Index above 2, located in a tropical

601    climatic area, and where no dams or hydraulic infrastructures are present. Also, out of our

602    locations, high latitude rivers with seasonal ice-cover tend to exhibit good performance.

603    Secondly, we performed a classification decision tree analysis, based on Random Forest, to obtain

604    the variable importance when classifying a site as good or poor. The output of this analysis showed

605    that mean observed discharge, climatic region, land cover and mean leaf area index (LAI) and

606    upstream catchment area and were the variables with higher importance, whereas river ice and

607    dam obtained the lowest importance. Both the individual and the combined classification analysis

608    of these local characteristics give us critical evidence of the relationship between the ground and

609    satellite discharge measurements and when it is expected to perform well. Furthermore, it provides

610    a guideline for future selection of measuring sites.

611    The locations with a very good performance will be selected for a potential future project where

612    satellite measure discharge could be calculated for longer periods and on a daily basis from the

613    remote sensing signal, analogous to the Dartmouth Flood Observatory method. This will represent

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

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

## 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.

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

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

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

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.

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

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

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

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

Disaster Charter, 2013. Space and Major Disasters. http://www.disasterscharter.org/, last accessed 1 September 2014.

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

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

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

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

Global Runoff Data Centre, The. River Discharge Time Series. 56068 Koblenz, Germany: Federal Institute of Hydrology (BfG). http://grdc.bafg.de/, last accessed 20 January, 2013

Golnaraghi M., J. Douris, J.-B. Migraine: Saving Lives Through Early Warning Systems and Emergency Preparedness, Risk Wise, Tudor Rose, pp 137–141, 2009.

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

Hirpa FA, Hopson TM, De Groeve T, Brakenridge GR, Gebremichael M, Restrepo PJ. Upstream satellite remote sensing for river discharge forecasting: Application to major rivers in South Asia. Remote Sens Environ, 131:140-51, 2013.

Hirsch R. M, Costa J. E.:U.S. Stream Flow Measurement and Data Dissemination Improve EOS, Transactions, American Geophysical Union. Vol. 85, No. 20, 18 May 2004, 197-203 pp, 2004.

Khan, S.I., Hong, Y., Vergara, H.J., Gourley, J.J., Robert Brakenridge, G., De Groeve, T., Flamig, Z.L., Policelli, F. & Yong, B.: Microwave satellite data for hydrologic modeling in ungaued basins, IEEE Geoscience and Remote Sensing Letters, vol. 9, no. 4, pp. 663-667, 2012.

Kugler, Z., and De Groeve, T.: The Global Flood Detection System, Office for Official Publications of the European Communities, Luxembourg, 2007.

Kugler, Z.: Remote sensing for natural hazard mitigation and climate change impact assessment, Quaterly Journal of the Hungarian Meteorological Service. Vol.116, No.1, January-March 2012, pp.21-38, 2012.

Kundzewicz, Z. W: Changes in Flood Risk in Europe, Wallingford: IAHS Press. 516 p. IAHS special publication; 10, 2012. United Nations: Report of the United Nations Conference on Sustainable, Development Rio de Janeiro, Brazil. 20–22 June 2012, A/CONF.216/16, 2012.

Le Coz, J., Renard, B., Bonnifait, L., Branger, F. & Le Boursicaud, R.: Combining hydraulic knowledge and uncertain gaugings in the estimation of hydrometric rating curves: A Bayesian approach, Journal of Hydrology, vol. 509, pp. 573-587, 2014.

Lehner, B., and Döll, P.: Development and validation of a global database of lakes, reservoirs and wetlands. Journal of Hydrology 296/1-4: 1-22, 2004. doi: 10.1016/j.jhydrol.2004.03.028

Lehner, B., Reidy Liermann, C., Revenga, C., Vörösmarty, C., Fekete, B., Crouzet, P., Döll, P., Endejan, M., Frenken, K., Magome, J., Nilsson, C., Robertson, J., Rödel, R., Sindorf, N., Wisser, D.: High resolution mapping of the world's reservoirs and dams for sustainable river flow management, Frontiers in Ecology and the Environment. Source: GWSP Digital Water Atlas. Map 81: GRanD Database (V1.0), 2008. Last access: 11/03/2014. http://atlas.gwsp.org/index.php?option=com_content&task=view&id=209&Itemid=1

Liaw, A. and Wiener, M. Classification and Regression by randomForest. R News 2(3), 18—22, 2002.

Moffitt, C.B., F. Hossain, R.F. Adler, K.K. Yilmaz, and H.F. Pierce. Validation of a TRMM-Based Global Flood Detection System in Bangladesh. International Journal of Applied Earth Observation and Geoinformation Volume 13, Issue 2, April 2011, Pages 165-177, DOI: 10.1016/j.jag.2010.11.003.

MunichRe, Munich Reinsurance: January 2014 press release, 2014.Münchener Rückversicherungs-Gesellschaft, Geo Risks Research, NatCatSERVICE http://www.munichre.com/en/media_relations/press_releases/2014/2014_01_07_press_rele ase.aspx, last access 20 January 2014

Nash, J.E. & Sutcliffe, J.V.: River flow forecasting through conceptual models part I - A discussion of principles Journal of Hydrology, vol. 10, no. 3, pp. 282-290, 1970.

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 behavior of random forest permutation-based variable importance measures under predictor correlation. BMC Bioinformatics, 11:110, 2010. doi: 10.1186/1471-2105-11-110

Palczewska, A., Palczewski, J., Robinson, R.M. and Neagu, D.: Interpreting random forest models using a feature contribution method, Proceedings of the 2013 IEEE 14th International Conference on Information Reuse and Integration, IEEE IRI 2013, pp. 112, 2013.

Pappenberger, F., Matgen, P., Beven, K.J., Henry, J.B., Pfister, L., de Fraipont, P. Influence of uncertain boundary conditions and model structure on flood inundation predictions. Adv. Water Resour. 29, 1430–1449, 2006. doi: 10.1016/j.advwatres.2005.11.012

Pappenberger, F., Dutra, E., Wetterhall, F. & Cloke, H.L.: Deriving global flood hazard maps of fluvial floods through a physical model cascade, Hydrology and Earth System Sciences, vol. 16, no. 11, pp. 4143-4156, 2012.

Peel, M. C., Finlayson, B. L., and McMahon, T. A.: Updated world map of the Köppen-Geiger climate classification, Hydrol. Earth Syst. Sci., 11, 1633-1644, doi: 10.5194/hess-11-1633-2007, 2007.

Pelletier, P.M.: Uncertainties in the single determination of river discharge: a literature review. Canadian Journal of Civil Engineering, 15:834–850, 1988.

Rosso, R. A linear approach to the influence of discharge measurement error on flood estimates. Hydrol. Sci. J. 30 (1), 137–149, 1998. doi: 10.1080/02626668509490975

Sandri, M. and Zuccolotto, P.. A bias correlation algorithm for the Gini variable importance measure in classification trees. Journal of Computational and Graphical Statistics, 17:611-628, 2008. [184], doi: 10.1198/106186008X344522

Schumann, Guy, Paul D. Bates, Matthew S. Horritt, Patrick Matgen, and Florian Pappenberger.: Progress in Integration of Remote Sensing–derived Flood Extent and Stage Data and Hydraulic Models. Reviews of Geophysics 47, RG4001, no. 4, 2009. doi: 10.1029/2008RG000274.

South African Water Affairs (DWA) database, http://www.dwa.gov.za/Hydrology/, last access: 10 July 2013.

Strobl, C., Malley, J. & Tutz, G.: An Introduction to Recursive Partitioning: Rationale, Application, and Characteristics of Classification and Regression Trees, Bagging, and Random Forests, Psychological methods, vol. 14, no. 4, pp. 323-348, 2009. doi: 10.1186/1471-2105-9-307

Thielen, J., Bartholmes, J., Ramos, M.-H., and de Roo, A.: The European Flood Alert System. Part 1: Concept and development, Hydrol. Earth Syst. Sci., 13, 125–140, doi: 10.5194/hess-13-125-2009, 9278, 2009.

Thiemig, V., Bisselink, B., Pappenberger, F., and Thielen, J.: A pan-African Flood Forecasting System, Hydrol. Earth Syst. Sci. Discuss., 11, 5559-5597, doi: 10.5194/hessd-11-5559-2014, 2014.

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

Tomkins, K.M.: Uncertainty in streamflow rating curves: Methods, controls and consequences. Hydrological Processes, 28(3), pp. 464-481, 2014.

UNISDR: Global Assessment Report: Revealing Risk, Redefining Development, Chapter 2.2. Global disaster risk trends, United Nations, printed in the UK, ISBN 978-92-1-132030-5, page 22-27, 2011.

UNOSAT, UNITAR Operational Satellite Applications Programme http://www.unitar.org/unosat/maps, last accessed 1 December 2013

Van Westen, C.J.: Remote sensing and GIS for natural hazards assessment and disaster risk management. In: Shroder, J. (Editor in Chief), Bishop, M.P. (Ed.), Treatise on Geomorphology. Academic Press, San Diego, CA, vol. 3, Remote Sensing and GIScience in Geomorphology, pp. 259–298, 2013.

Yamazaki, D., O'Loughlin, F., Trigg, M.A., Miller, Z.F., Pavelsky, T.M. & Bates, P.D. 2014, "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.

Zaraj, Z., Zambrano-Bigiarini, M., Salamon, P. Burek, P., Gentile, A., and Bianchi, A.: Calibration of the LISFLOOD hydrological model for Europe. Calibration Round 2013JRC Technical Report, European Commission, Joint Research Centre. Ispra, Italy, 2013.

Zhang, Y., Hong, Y., Wang, X., Gourley, J.J., Gao, J., Vergara, H.J. and Yong, B.: Assimilation of passive microwave streamflow signals for improving flood forecasting: A first study in Cubango River Basin, Africa. IEEE Journal of Selected Topics in Applied Earth 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.[1]Stations 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 Catchment[1] | Upstream catchment areas (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. [2]Types 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 |

6

1    **Table A 1.** Studied land cover types from GlobCover (2009) aggregated into broader categorical

2    classes by type and vegetation density.

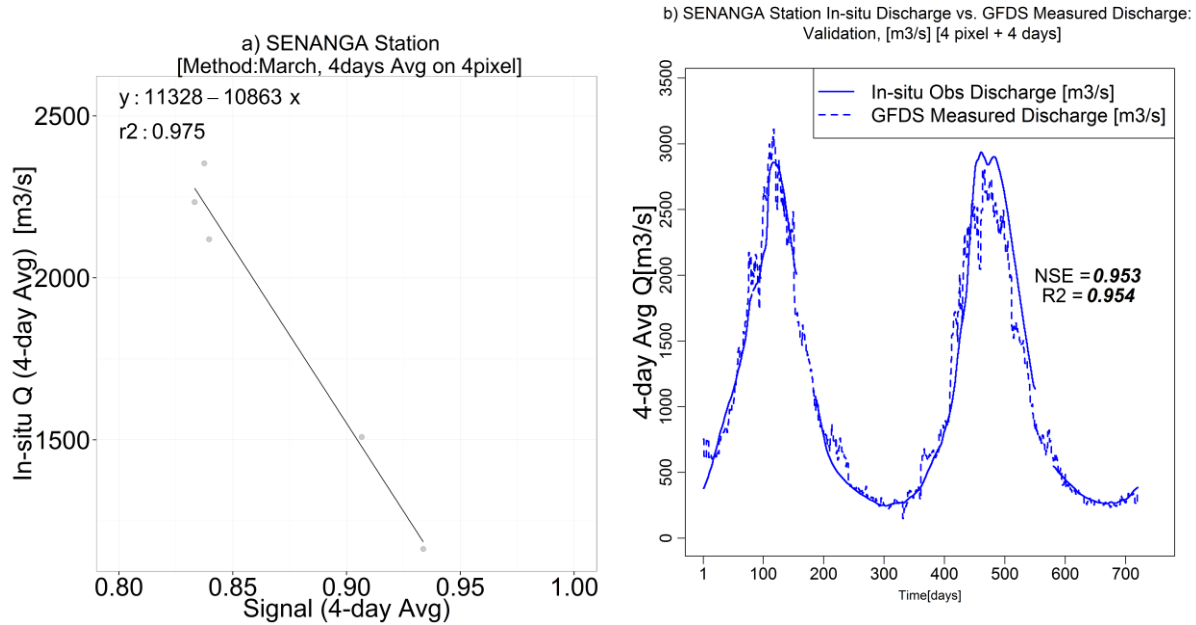| Label | Aggregated classes |
| --- | --- |
| Rainfed croplands | Rainfed croplands |
| Sparse (<15%) vegetation | Sparse vegetation |
| Closed to open (>15%) broadleaved evergreen or semi-deciduous forest (>5m) | Closed to open forest |
| Closed to open (>15%) mixed broadleaved and needleleaved forest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

3

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.

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

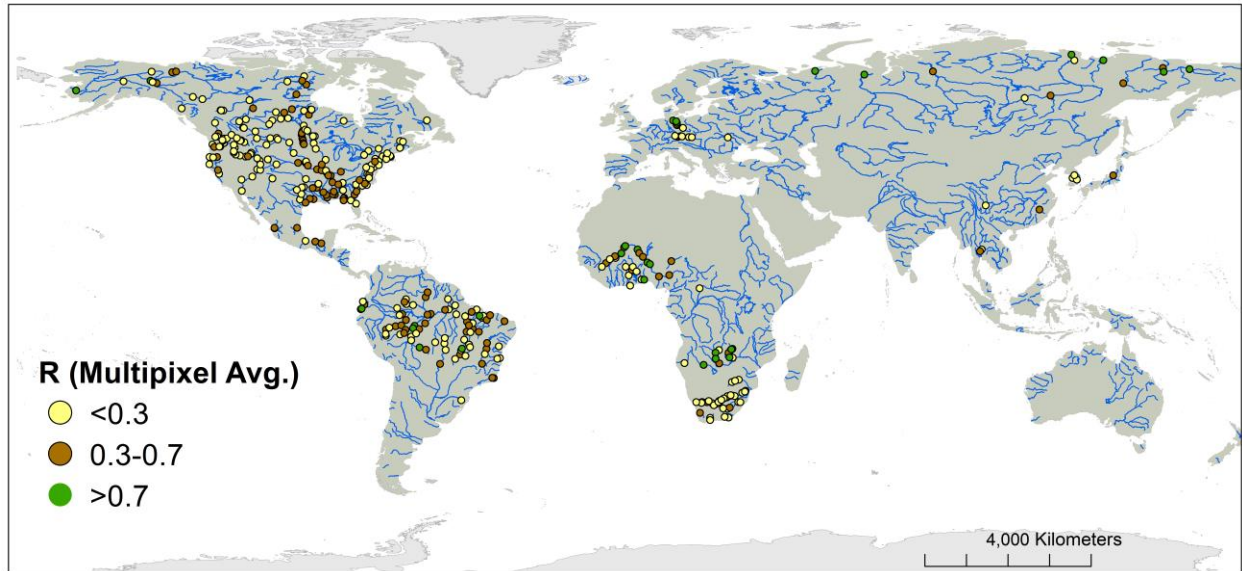**Figure 3.** a) scatterplot for the Senanga station (Long 23.25, Lat. -16.116) in the Zambezi River (Africa). Monthly mean for March from 1998 up to 2002. b) Validation hydrograph for 2003-2004 and skill scores for Senanga. The (monthly) rating equations were used to calibrate the signal into discharge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

1

**Figure 4**. Location of stations and R skill score between in situ observed discharge and satellite
signal (4 days and 4 pixels average). Globally, 169 sites have R>0.3, of which 42 have R>0.5.

1

**Figure 5.** Nash-Sutcliffe efficiency of the validation (n= 332 stations). Globally, 154 stations have
NSE>0 of which 80 stations have NSE> 0.50.
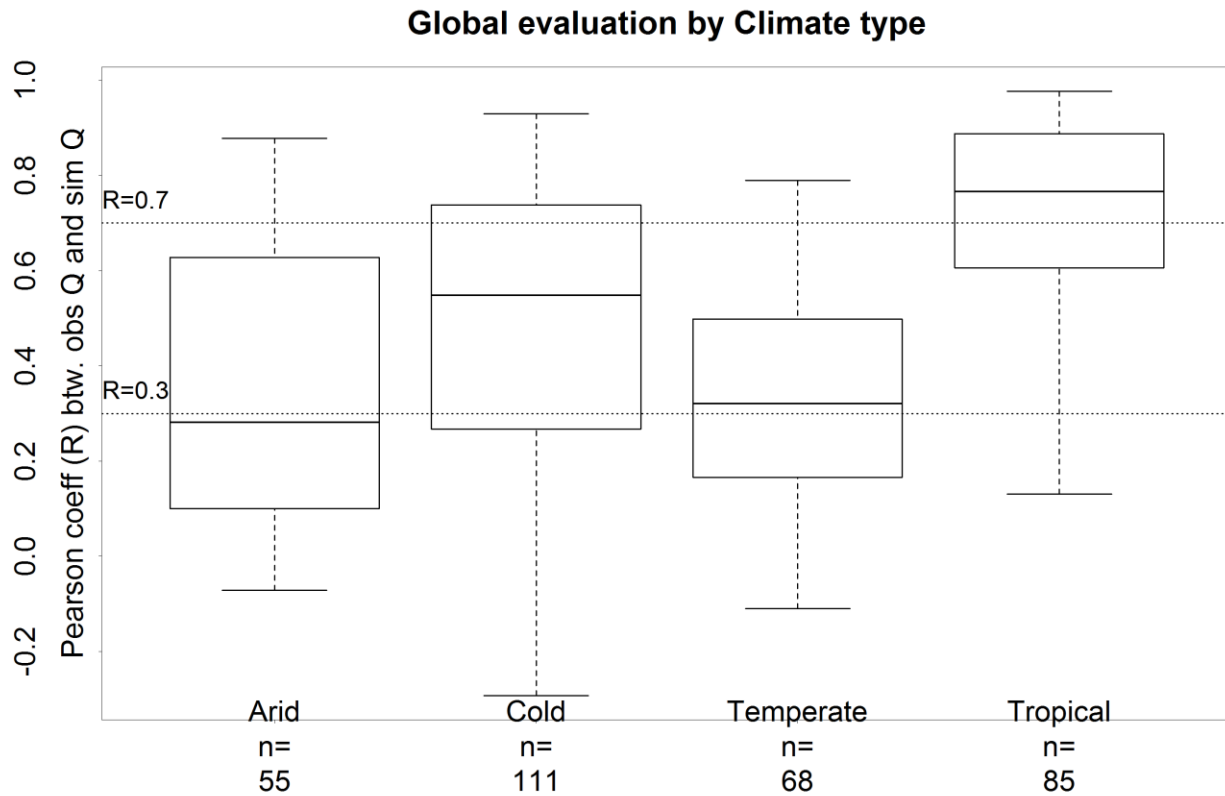
**Figure 6**. a) relationship between R obtained from the validation of satellite measured discharge and the maximum river width for each location; b) relationship between the same R score and the presence of significant floodplains, flooded forest and wetlands Horizontal dotted line shows the R=0.3 and R=0.7 threshold, the vertical line is the river width equal to 1km.

1

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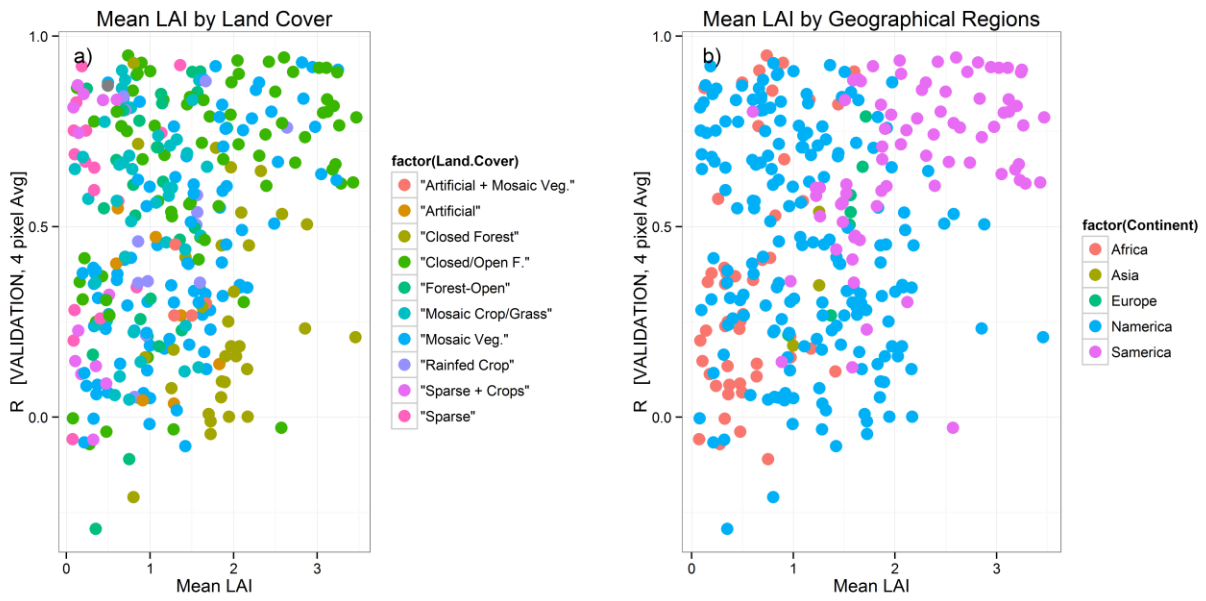
**Global evaluation by Land Cover type**

1

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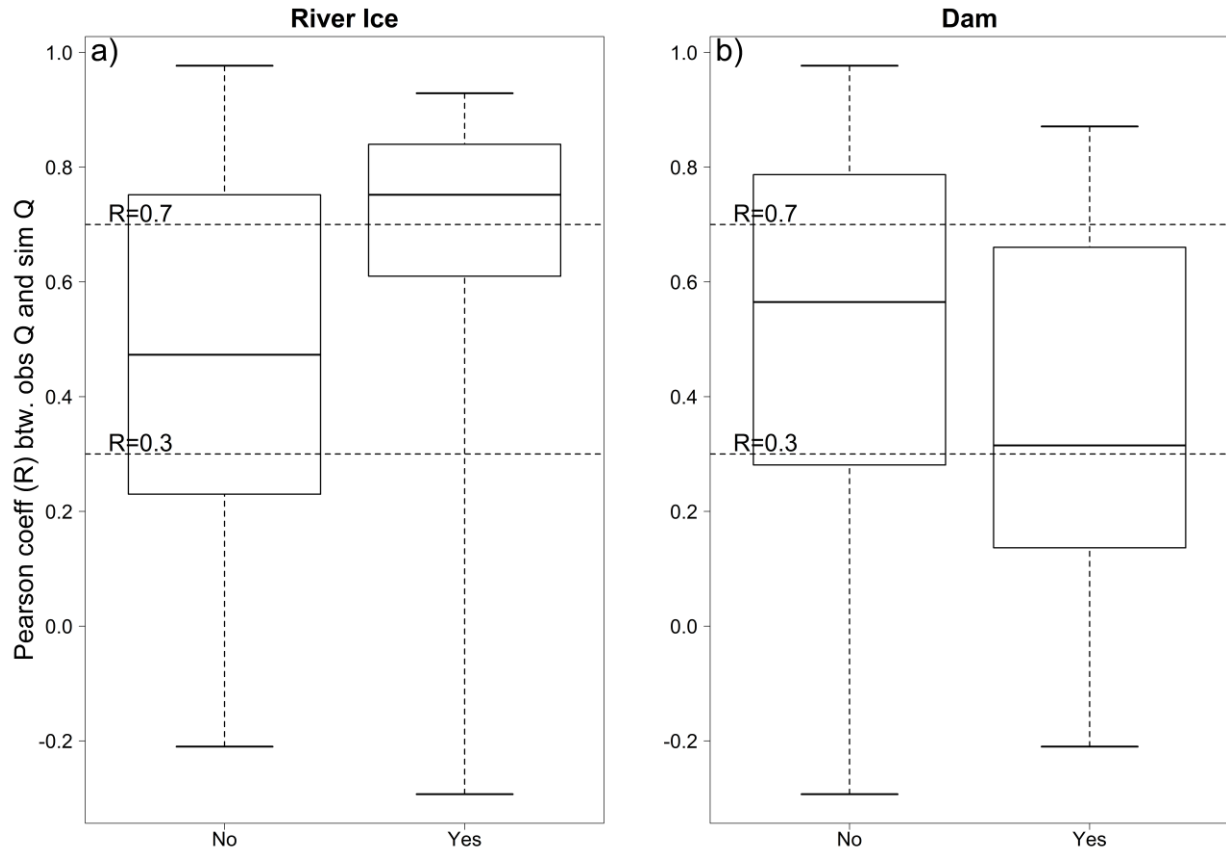
**Global evaluation by Climate type**

1

**Figure 9.** Global evaluation of the R score obtained during the validation and its classification –
only main types-by the Köppen-Geiger climate area (Peel et al. 2007). Note that polar climate was
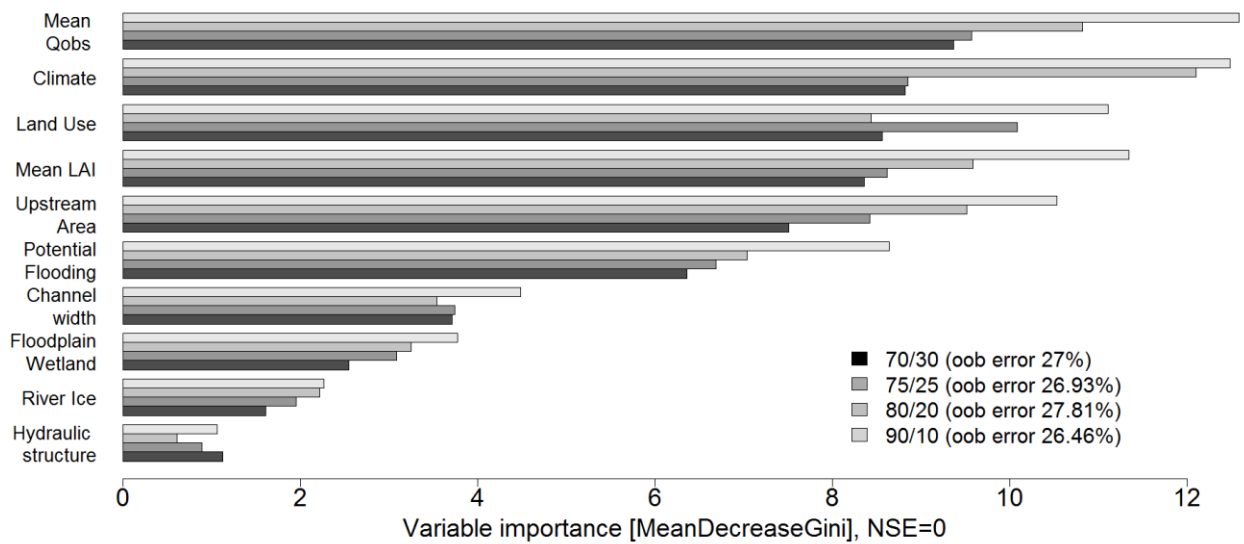excluded from this analysis as only three stations felt into this category.

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

1

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

1

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