Dear Editor,

We would like to thank the two anonymous referees for their valuable comments and suggestions to enhance the scientific quality of our contribution. We deeply appreciate their enthusiastic comments to the manuscript. In the following, we have addressed the reviewers' comments and provided evidence (red text) of the corresponding revisions to the manuscript.

Kind regards on behalf of all authors,

Maurizio Mazzoleni

Anonymous Referee #1:

This contribution deals with the question how crowdsourced observations could be utilized in hydrological modelling. The approach used here is to use such observations to update models used for forecasting runoff. The challenge is that the observations might come at irregular times and with varying accuracy. This is an interesting and timely issue and I was excited when I started reading the manuscript. In the end, however, I have to admit, I was not fully convinced, and feel that a major revision is needed.

RC: My major is the selection of catchments/models and limited event data being used here. Two catchments were chosen (the selections seem a bit random, but ok) and two different models were used for the two catchments. The latter seems to make little sense, as it makes results hardly comparable. It should be noted that also issues like calibration largely varied: in the Brue case, calibration was based on one event only (p11381,l12), whereas in the Bacchiglione catchment ten years were used for calibration (p11383,l20). In the first case some form of effective precipitation must have been used, as only the so called direct runoff is simulated (but it is unclear how this was determined), whereas in the second case the entire runoff has been simulated.

AC: Thank you for this comment. We realise we were not fully clear in the manuscript. The mentioned important points are now clarified in the revised manuscript. Ideally, a new method has to be tested on a variety of catchments and models, but in the presented study we were constrained by the available funds, the projects limitations and the data availability.

We have chosen these two case studies because we wanted to assess if the results we have obtained are valid for areas with different topographical and hydrometeorological features and represented by two different hydrological models.

First, the Bacchiglione basin is one the official cases studies of the WeSenseIt project which is funding this research. We had the model available from AAWA, the water authority (Ferri et al. (2012)). The developed methodology has been tested using synthetic observations, and the intention is to then apply it on the existing early warning system on the Bacchiglione Basin, previously designed and implemented by Ferri et al. (2012).

The second case, the Brue catchment, is an experimental catchment and is used in many studies, including us (Mazzoleni et al. (2015)). In addition, we specified in section 3.1 that, in case of the Brue catchment:

"The choice of the model is based on previous studies performed on the Brue catchment in case of assimilation of streamflow observations from dynamic sensors (Mazzoleni et al., 2015)"

We agree the cases and models are indeed different, but the presented study demonstrated that the results obtained, Figure 11 and 13, are very similar in terms of model behaviour assimilating asynchronous observations. This does not allow for a universal generalisation (that the methods works in all cases), and we are not claiming this. In the present version we clearly state that:

...additional analyses on different case studies and larger time series of flood event should be carried out in order to provide more general conclusions.

These additional clarifications have been added as in the introduction:

"The Brue catchment is considered because of the availability of precipitation and streamflow data, while the Bacchiglione river is one of the official case studies of the WeSenseIt Project (Huwald et al., 2013), which is funding this research"

and conclusion:

"We agree the cases and models are indeed different, but the presented study demonstrated that the results obtained are very similar in terms of model behaviour assimilating asynchronous observations"

On the data sets used for calibration, we understand the reviewer concern regarding the short time series used to calibrate the lumped model in the Brue catchment. For this reason, an additional calibration is performed considering the period from 23-10-1994 to 17-03-1995, affected by a long series of flood events, to demonstrate that model parameters estimated in this new calibration do not differ from the ones previously estimated. In both calibration we used the correlation R as objective function of the optimization method. The result of the calibration performed in case of longer time series is showed in Figure 1. Table 1 shows the calibrated values of the parameter c_k and n obtained using two different time series.

Table 1 Parameters values in case of two different time series used to calibrate the model

| | <i>n</i> (-) | c (-) |
|--------------------|--------------|-------|
| Short time series | 5 | 0.035 |
| Longer time series | 4 | 0.026 |

As it can be seen from Table 1, similar results are obtained in both cases for short and long time series. In particular, no difference is visible in the flood event 1 and 2 when applied the two different calibrated parameter sets. For this reason, the text in section 3.1 of the revised manuscript have been changed as:

"The model calibration is performed maximizing the correlation between the simulated and observed value of discharge, at the outlet point of the Brue catchment, during the flood events occurred from the 23-10-1994 to 17-03-1995. The results of such calibration provided a value of the parameters n and c_k equal to 4 and 0.026 respectively"



Figure 1. Calibration of Brue model in case of longer time series.

In case of the Bacchiglione catchment, as already stated, the model was already calibrated by AAWA in order to be used in the early warning system for the Bacchiglione River.

In order to clarify the input and output of the conceptual model used in the Brue catchment, the following text in the revised manuscript has been added:

"where I is the model forcing (in this case direct runoff), n (number of storage elements) and k (storage capacity) are the two parameters of the model and Q is the model output (streamflow)."

Direct runoff is estimated using the approach used in Mazzoleni et al (2015) as referred in the text. On the other hand, in the Bacchiglione model the input is the time series of precipitation. **RC:** The discussion of the models, especially the second one, largely ignores recent findings on runoff generation processes. For instance, the statement on residence time (P11383,110) should be reformulated with the recent paper of Beven and McDonnell in mind. In the end, for this study the physical correctness of the models is probably less important, but I still find the uncritical description of the models with their partly unrealistic assumptions a bit troublesome.

AC: Indeed, descriptions and limitations of the models used had to be made more explicit. Following this reviewer comment we revised the model description related to the runoff generation process. The following text has been added to the revised version of the manuscript (sec. 3.2):

"In particular, in case of Qsur the value of the parameter k, which is a function of the residence time in the catchment slopes, is estimated relating the slopes velocity of the surface runoff to the average slopes length L. However, one of difficulties involved is proper estimation of the surface velocity, which should be calculated for each flood event (Rinaldo and Rodriguez-Iturbe, 1996). According to Rodríguez-Iturbe et al. (1982), such velocity is a function of the effective rainfall intensity and event duration. In this study, the estimate of the surface velocity is performed using the relation between velocity and intensity of rainfall excess proposed in Kumar et al. (2002). In this way it is possible to estimate the average time travel and the consequent parameter k. However, such formulation is applied in a lumped way for a given sub-catchment. As reported in McDonnell and Beven (2014) more reliable and distributed models should be used to reproduce the spatial variability of the residence times within the catchment over the time. That is why, in the advanced version of the model implemented by AAWA, in each subcatchment the runoff propagation is carried out according to the geomorphological theory of the hydrologic response. In such model, the overall catchment travel time distributions is considered as nested convolutions of statistically independent travel time distributions along sequentially connected, and objectively identified, smaller sub-catchments. The parameter k assumes different values for each time step as the rainfall changes. In fact, the variability of residence time is considered according to Rodríguez-Iturbe et al. (1982) by assuming the surface velocity as a function of the effective rainfall intensity (Kumar et al., 2002). Anyway, the correct estimation of the residence time should be derived considering the latest findings reported in McDonnell and Beven (2014). In case of Qsub and Qg the value of k is calibrated comparing the observed and simulated discharge at Vicenza as previously described."

The goal of this study is to assess the effect of assimilating crowdsourcing observations coming at irregular time steps within an existing early warning system to then apply such methodology in real-life. For this reason, we wanted to keep the Bacchiglione model used in this study as close as possible to the one implemented in the AMICO platform used by AAWA.

RC: Most importantly, however, I find the small number of tested events problematic (2 in Brue, 1 in Bacchiglione). Obviously the results depend largely on the characteristics of the event and the quality of the precipitation data. I am afraid that this extremely small number of events makes results rather 'random'. Honestly, I find it therefore difficult to see what this study contributes beyond that the additional information improves simulation somewhat (which one would have expected anyway).

AC: We agree with the reviewer that the small number of flood events might be problematic to draw general conclusions. We have revised the conclusions and added the following (a limitation of this study):

"...additional analyses on different case studies and the longer time series of flood events should be carried out in order to draw more general conclusions about assimilation of the crowdsourced observations and their value in different types of catchments and model setups".

and formulated the corresponding recommendation in the end of sec. 7.

At the same time we would like to stress that we found similar trends in the dependency of Nash index on the number of observations - in both flood events in Brue (Figure 2, Figure 3 and Figure 4) and Bacchiglione. It is now stated in Sec. 6.1 clearer. In fact, after a threshold number of observations, NSE asymptotically approaches a certain value, in both flood events, meaning that no improvement is achieved with additional observations (Figure 2). However, the only difference is that threshold number of observations and asymptotic NSE values are different because model performances can change according to the considered flood events. Example of these considerations are showed in the following two figures



Figure 2. Model improvement during flood event 1 and 2, in case of different lead times, assimilating streamflow observations according to scenario 1



Figure 3. Dependency of μ (NSE) and σ (NSE) on the number of observations, for the scenarios 2, 3, 4, 5, 6, 7, 8 and 9 in case of flood event 1.



Figure 4. Dependency of μ (NSE) and σ (NSE) on the number of observations, for the scenarios 2, 3, 4, 5, 6, 7, 8 and 9 in case of flood event 2.

This is one of the first studies that address the issue of including crowdsourcing observations coming at asynchronous instant within hydrological modelling. We believe that this study demonstrated that assimilation of crowdsourcing observations can improve hydrological modelling requesting a limited amount of observations, depending on the flood event, even if such observations have variable uncertainty and are coming at random moment.

Our ongoing research is actually looking at a higher number of flood events. After discussions between the authors, we concluded that the current manuscript is already too long to include these new analyses, which will require a number of extra figures.

RC: The more interesting questions of how big the improvement is, how many observations are needed, at which accuracy, \ldots all are too heavily influenced by the choice of the one or two event(s) to be of a more general value.

AC: We agree with reviewer comment on the importance of knowing the amount of needed observations and their level of accuracy. However, model performances vary according to the type of flood event as showed in this study. Consequently, improvement should be related to the particular flood event. It is not possible to define a priori number of observations needed to improve model. In fact, considering figure 1 reported in this reviewer response, in case of flood event 1, we needed only 5 observations to reach NSE equal to 0.9, while during flood event 2, we needed 15 observations to reach the same NSE value. That is why the number of observations necessary to achieve a good model performance might vary according to the flood event. A possible improvement of our method could be a pre-filtering module aimed to select only observations having good accuracy while discarding the ones with low accuracy. However, in this study we aim to assimilate all type of observations without giving recommendations on the minimum acceptable accuracy value of crowdsourced observations. In this way, we demonstrated that high number of observations are included in the revised manuscript.

In sec. 6.1 the following text has been added:

"In both flood events we found similar trends in the dependency of Nash index on the number of observations. However, it is not possible to define a priori number of observations needed to improve model. In fact, after a threshold number of observations (five for flood event 1 and fifteen for flood event 2), NSE asymptotically approaches to a certain value meaning that no improvement is achieved with additional observations."

Minor issues:

RC: The language could be improved, there are several small language mistakes, which make reading more difficult.

AC: Language has been improved in the revised version of the manuscript.

RC: *The graphs could be improved, they are in general quite hard to read (and not be too 'nice' to be honest)*

AC: Indeed, legends on some Figures were not clear. We have improved the quality of Figures 2, 3, 4, 6 and 10. Special attention was given to revising Figure 3 with the scheme of the hydrological model used in this study.

RC: Use mathematically correct terms in your equations! ET(Eq.2), for instance, is not correct as it strictly mathematically means E times T (note that you actually use this in the directly following equation, where CS(t) actually means C times S(t))

AC: We modified the equations of the manuscript according to reviewer comment

RC: Please separate results and discussion; this would make reading the text so much easier.

AC: We have seriously considered this suggestion. Indeed, it is a possibility – but still we would like to suggest that we would follow the initial logic of having results and discussion integrated in one section. This is in a way subjective decision but we checked with a colleague and all came to a conclusion that in this paper the results are presented clearer if each of them is immediately followed by interpretation and corresponding discussion.

RC: Be careful with your references, some are missing in the reference list (e.g. Krouse) other are misspelled (e.g. Bergström)

AC: We revised the reference list and corrected the misspelled and added the missing references.

RC: Reference to recent work could be improved. The work could be better linked to recent work on the value of (limited) data in hydrological modeling. Also, the recent review on citizen science in hydrology by Buytaert et al. (2014) should be referred to.

AC: We thank the reviewer for the useful comment. This paper deals only with citizen data used for hydrological modelling, and we do not study the value of information in the traditional sensor data. We know of attempts in France and Italy to use the videos of the flood events to estimate the water velocity and to use this information to re-calibrate the river models, however could not find the published papers on this subject. The only reference we could find is EGU abstract by Candela et al. (2013) where citizen data (video) was used to validate the urban flood model in a post-event analysis comparing water depth or flow velocity. To the best of our knowledge, there are no studies that would quantitatively analyse the usefulness of using citizens' data for (real-time) DA in flood modelling.

We do not want to extend the scope of the paper and to deal only with assimilation of crowdsourced observations. We have also followed an advice of reviewer #2 to shorten the Introduction. The reference to Buytaert et al (2014) is very useful indeed (however it does not consider the examples of DA of citizen data for hydrology). We added this references, and the reference to the CitiSense project (this project belongs to the pool of five Citizen Observatories EU projects, together with WeSenseIt which funds this study).

Candela, A., Naso, S. And Aronica, G.: On the use of innovative post-event data for reducing uncertainty in calibrating flood propagation models, EGU General Assembly, Vienna, Austria, 7-12 April, 2013

References:

- McDonnell, J. J., and K. Beven (2014), Debates—The future of hydrological sciences: A (common) path forward? A call to action aimed at understanding velocities, celerities and residence time distributions of the headwater hydrograph, Water Resour. Res., 50, 5342–5350, doi:10.1002/2013WR015141
- Buytaert, W., Zulkafli, Z., Grainger, S., Acosta, L., Bastiaensen, J., Bhusal, J., Clark, J., Dewulf, A., Foggin, M., Hannah, D.M., Hergarten, C., Isaeva, A., Pandey, B., Paudel, D., Sharma, K., Steenhuis, T., Tilahun, S., Van Hecken, G., Zhumanova, M., 2014. Citizen science in hydrology and water resources:opportunities for knowledge generation, ecosystem service management, and sustainable development. Front. Earth Sci. doi:10.3389/feart.2014.00026

Anonymous Referee #2:

Thank you for the opportunity to review the article "Can assimilation of crowdsourced streamflow observations in hydrological modelling improve flood prediction?" (hess- 2015-415). This article presents an evaluation of methods for improving the accuracy of hydrologic models by incorporating crowdsource (social sensors) data. This is an interesting idea and the first paper on the topic that I have read. The opportunity to get the public to engage in extreme-events using technology they are already familiar with is exciting and will likely be a great success. I think the paper is generally written well and accurately presents the methods and results and that the discussion and conclusions are reasonable. That said, I have included a few comments/suggestions/questions for the authors to consider. I have not provided an editorial review, though I do believe the paper should have a thorough editorial review prior acceptance. There are several instances with subject / verb agreement, some words are unnecessarily plural, and acronyms that do not appear to have definition (DA for example). Additionally, figures need to be checked to make sure they include relevant information included in the text (for example, include "setting A" on figure 15 or describing (a) and (b) on figure 13).

RC: Are there any methods currently in use to quantify the accuracy of crowdsource (CS) data? This is particularly important given that the methods you for including crowdsourced data are workable. I think you mention briefly about assessing accuracy of actual social sensors. Please expand on this in terms of current ideas, particularly ideas that would assess accuracy in an objective manner

AC: We thank the reviewer for this valuable comment. Following his suggestion, we included the following additional information about methods used to assess quality of CS data, in the introduction:

"According to Bordogna et al. (2014) and Tulloch and Szabo (2012), quality control mechanisms should consider contextual conditions to deduce indicators about reliability (expertise level), credibility (volunteer group) and performance of volunteers such as accuracy, completeness and precision level. Bird et al. (2014) addressed the issue of data quality in conservation ecology by means of new statistical tools to assess random error and bias in such observations. Cortes et al. (2014) evaluated data quality by distinguishing the in-situ data collected between a volunteer and a technician and comparing the most frequent value reported at a given location. They also gave some range of precision according to the rating scales. With in-situ exercises, it might be possible to have an indication of the reliability of data collected (expertise level). However, this indication does not necessarily lead to a conclusion of high, medium or low accuracy every time a streamflow observation

of a contributor is received. In addition, such approach is not enough at operational level to define accuracy in data quality. In fact, every time a crowdsourced observation is received in real-time, the reliability and accuracy of observations should be identified. To do so, one possible approach could be to filter out the measurements following a geographic approach which defines semantic rules governing what can occur at a given location (e.g. Vandecasteele and Devillers, 2013). Another approach could be to compare measurements collected within a pre-defined time-window in order to calculate the most frequent value, the mean and the standard deviation."

RC: Please consider restructuring the Introduction. While the Introduction is very informative, it is quite long and digresses into a discussion of sensor technologies, issues of quality control, other CS networks, oceanographic models, and assimilation of asynchronous observations among other things. The paper is supposed to be about assimilation of CS data assimilation. The Introduction should go directly to this point. As written, the introduction of the topic and explanation of the objectives are separated by a considerable amount of material. Please shorten the Introduction to clearly present the topic, current understanding of how to include CS data, gaps in that understanding, and what you propose to do to fill that knowledge gap. The other information should be retained, but put into a different sections ("Background", "Existing CS Networks"). I personally find the material on existing CS networks very interesting and would like to see that information discussed a bit more.

AC: Following reviewer's suggestion we shortened the introduction and focused on assimilation of CS observation, providing also some details about the past and ongoing projects in which CS are used to improve models predictions. In particular, we firstly defined the necessity of improving the model introducing the concept of model updating and DA. Secondly, we described the need of CS observations and some CS projects are illustrated. We focused on two main characteristics of the crowdsourced observations: a) data quality and b) variable life span (asynchronous observations). Thirdly, data quality issues and method used to deal with this problem are described (following previous reviewer's comment). Finally, we described existing methods used to assimilated asynchronous observations in hydrology and other water related models. The text related to the assimilation of distributed hydrological observations has been removed. We believe that in the present form the introduction is more readable and objectives of this paper are clearer. We also provided additional details about methods to assess observational uncertainty as proposed by reviewer in a previous comment. The new version of the introduction is included in the revised manuscript.

RC: Is the discussion about oceanographic studies / models needed? It was not clear to me what that material added to the paper. If it is needed, please make it more clear what the connection is? Is it technology of oceanographic models that can be used in your process of including CS data into models?

AC: DA in oceanographic models is in the paper since oceanographic observations are commonly collected at not pre-determined, or asynchronous, times – and this has relevance for the paper. Indeed, the DA technologies used in oceanography (continuous (variational DA)) could have been used also for hydrology, but they require building adjoint models and this limits their use in case of using real complex hydrological models. (The reasons of using KF instead are given in Introduction.)

RC: Why is the MIKE11 model presented as the model for representing flood propagation on the main channel in the Bacchiglione basin? Immediately after stating that the MIKE11 model was used, it appears that it was replaced by the Muskingum - Cunge model. Maybe they were used to represent two different processes in this basin? Obvioulsy, this was not clear. If you used the M-C model, then why even bother with the MIKE11 part of the discussion? Please reconsider your wording to make it clear. If both were used, please explain the role of each.

AC: We thank the reviewer for pointing out this aspect of our study, which perhaps was not clearly explained. MIKE11 model was originally used by AAWA, the water authority, within their early warning system on the Bacchiglione basin. However, in this study, in order to reduce the computational time of the simulations and since main part of the uncertainty sources come from the hydrological model, MIKE11 was replaced with a Muskingum-Cunge model. We mentioned this point in section 3.2:

"In the early warning system implemented by AAWA in the Bacchiglione catchment, the flood propagation along the main river channel is represented by one-dimensional hydrodynamic model, MIKE 11 (DHI, 2005). This model solves the Saint-Venant equations in case of unsteady flow based on an implicit finite difference scheme proposed by Abbott and Ionescu (1967). However, in order to reduce the computational time required by the analysis performed in this study MIKE11 is replaced by a hydrological routing Muskingum-Cunge model (see, e.g. Todini 2007), considering river cross-sections as rectangular for the estimation of hydraulic radius, wave celerity and the other hydraulic variables"

Due to limitation in the number of figures we do not present a comparison between MIKE11 and MC routing. However, we are currently working on a study in which we demonstrate that these two methods show similar results in terms of estimated discharge at Ponte degli Angeli.

RC: Increases in model accuracy due to assimilating CS observations needs to be presented in different ways. I understand the value in evaluating model accuracy and improvements in accuracy in terms of NSE. Several times in the paper, the value of including these CS observations is couched in terms of increased accuracy of flood peak magnitudes and timing. Discussing this increased accuracy in terms of NSE only is not all that informative. Statistics such as NSE only speak to

overall model accuracy, not to real increases/decreases in prediction error. Please include discussion about percent change in flood peak prediction (in text and/or table) for a few of the peaks in your evaluation period.

AC: Indeed, the averages statistics like NSE may not correctly present the model performance gains during floods, so the other error metrics which reflect flood-time performance more explicitly can be used. As suggested by the reviewer we have carried out additional analyses to assess the change in flood peak prediction considering 3 peaks occurred during flood event 2 (see Figure 3), in the Brue catchment.



Figure 1. Indication of the 3 flood peak occurred during flood event 2 in Brue catchment

Error in the flood peak timing and intensity is estimated using Err_t and Err_l equal to:

$$Err_t = t_P^o - t_P^S \,. \tag{1}$$

$$Err_{I} = \frac{Q_{P}^{o} - Q_{P}^{S}}{Q_{P}^{o}}.$$
(2)

Where t_p^o and t_p^s are the observed and simulated peak time (hours), while Q_p^o and Q_p^s are the observed and simulated peak intensity (m³/s). From the results in Figure 2 and 3, considering 12-hours lead time, it can be observed that, overall, errors reduction is achieved for increasing number of observations within 1 hour. In particular, assimilation of CS observations has more influence in the reduction of the peak intensity rather than peak timing. In fact, as Figure 4 shows, e.g. in case of peak 1, a small reduction of Err_t is obtained even increasing the number of CS observations. In fact, in all the 3 considered peaks, maximum reduction in Err_t is around 1 hour. On the other hand, higher error reduction is achieved if we considered the peak intensity rather than its timing. In particular. Smaller Err_t error values are obtained in case of scenario 1, while scenario 5 is the one that shows the lowest improvement in terms of peak prediction. This can be related to the random



moment and accuracy of the CS observation in such scenario. Similar results are obtained in case of scenario 6 and 9.

Figure 2. Representation of Err_t as function of the number of CS observations for 3 different peaks in case of scenarios from 1 to 9



Figure 3. Representation of *Err*₁ as function of the number of CS observations for 3 different peaks in case of scenarios from 1 to 9

These conclusions are very similar to the ones obtained analysing only NSE as model performance measures. This can be related to the linear nature of the model and the consequent DA approach used in this work. Due to the already high number of figures included in the revised version of manuscript, we have prepared an additional table (see below) indicating the percentage of error Err_t and Err_t reduction for each scenario changing from the assimilation of 1 to 20 observations. We leave the decision to the Editor whether to add or not these latest results and Figures/Table.

| | Err _t | | | Err _I | | |
|----------|------------------|----------|----------|------------------|----------|----------|
| Scenario | peak1 | peak2 | peak3 | peak1 | peak2 | peak3 |
| 1 | 0 | 0 | 0 | 0.588007 | 0.990132 | 0.545782 |
| 2 | 0.02 | 0 | -0.01 | 0.571863 | 0.97161 | 0.535791 |
| 3 | 0.015 | 0.069767 | 0.309524 | 0.529608 | 0.967312 | 0.480587 |
| 4 | 0 | 0 | -0.00337 | 0.564742 | 0.897512 | 0.535304 |
| 5 | 0.004975 | 0.186235 | 0.029126 | 0.486836 | 0.855197 | 0.443373 |
| 6 | 0.07 | 0 | -0.02 | 0.578038 | 0.969569 | 0.537394 |
| 7 | 0.05 | 0.052133 | 0.313333 | 0.528956 | 0.96833 | 0.481274 |
| 8 | -0.03093 | 0 | 0 | 0.560649 | 0.896826 | 0.535814 |
| 9 | 0.004975 | 0.177419 | 0.023333 | 0.489236 | 0.858807 | 0.441452 |

Table 1. Percentage of error *Err_t* and *Err_t* reduction