

1 **Reply to Anonymous Referee #1**

2 Note: We include our replies to the referee's comments directly at the respective points in the
3 text. Referee comments are kept in italics and our replies are kept in normal font style.

4 **1.1 General Comments**

5 *The paper addresses the problem of the credibility of grey literature on floods (which*
6 *presently hampers its implementation by researchers) and proposes a quality assessment*
7 *framework (QAF) for its evaluation. The point under investigation is a relevant technical*
8 *question (within the scope of NHESS) which has received little attention in the past, and*
9 *whose proper treatment could improve present capacity of analysts of understanding flood*
10 *risk.*

11 *The paper adapts existing tools, from other disciplines, to the problem under investigation*
12 *proposing, this way, an innovative method up to international standards. In this regards,*
13 *proper credit is given to previous work and authors' contribution is clear; reference is*
14 *appropriate and fully accessible by fellow scientists. In general, implemented data and*
15 *methods are clearly described but for statistical tools as better discussed in the following*
16 *section. The title and the abstract are appropriate; the presentation is well structured but the*
17 *paper is still too long. Figures and tables are adequate. The technical English language is*
18 *fluent and precise. Results and conclusion are significant but too much specifically related to*
19 *the case study. According to this, I suggest some revisions before paper's publication. In the*
20 *following general and specific comments are supplied.*

21 **1.2 Major criticisms**

22 *1) The "squared weighting scheme" implemented for the kappa test is not clear (section 2.3).*
23 *This makes difficult also the understanding of kappa test results (section 3.1). Both sections*
24 *should be re-written and made clearer.*

25 The squared weighing is a commonly adopted weighting scheme for ordinal scaled data. In
26 these data the neighbourhood of classes plays an important role as items that are allocated in
27 neighbouring classes adhere to more similarity than when they are allocated to opposite ends
28 of the scale. The weighting scheme considers these near-by characteristics and puts more
29 weight on neighbouring classes.

1 We agree that the current section on the kappa test is overly technical in its presentation and
2 will rewrite this section to clarify the relation to the data at hand. Accordingly we will revisit
3 the presentation of the results.

4 *2) Section 3.4: this section does not aim at demonstrating the applicability of the QAF, as*
5 *stated at the beginning of the section (this was already done in previous sections); rather, the*
6 *objective is to highlight how available reports and related information are actually useful to*
7 *understand/answer a specific technical question, if they are jointly used. The section should*
8 *then be re-written according to this perspective. Moreover, it should be put into light which is*
9 *the “weight” of information coming from reports with different quality in shaping the overall*
10 *information (credibility).*

11 The referee is perfectly right in this observation. The application highlights the potential of
12 combining information from many reports in order to understand a particular flood event. The
13 quality (overall and in the dimensions) of the reports is used to judge their applicability for the
14 task. However, it does not yet provide a framework for information expansion that includes
15 defined weights. This is subject for further research. We will consider this and rewrite this
16 section accordingly, i.e. not allude to the section as a demonstration of the QAFs applicability
17 rather than an illustrative example to highlight the potential of flood event documentation for
18 understanding a specific flood event.

19 In the concluding section (pg. 176, lines 14-23) we discuss the next steps needed to develop a
20 framework that formalizes the combination of information from many sources (i.e. combining
21 quality labelled information from event reports with model or data based analysis) and under
22 consideration of the uncertainties attached to each information. We will rewrite this paragraph
23 to make this clearer.

24 *3) Conclusions are too much related on German reports and their quality: this was already*
25 *(extensively) discussed in previous sections. Conclusions should be more generic, discussing*
26 *how the QAF can be implemented in research, with which improvements and consequences.*

27 We will consider the referees comment and rework the concluding chapter. I.e. we will
28 shorten the chapter substantially limiting the conclusions to only the most important findings
29 with a German specific notation and rather add more generic aspects. These are:

- 1 - the use of QAF for providing the basis for a better ad-hoc and post event analysis. What are
2 the critical factors that need to be analysed in the course of an event and what are critical
3 considerations that need to be taken in the design of a report
- 4 - the contribution/addition of this study to event databases/catalogues and the improvement to
5 an event set of floods by providing additional structured and quality labelled information.
- 6 - recommendations to report producers (better reporting can help improving capacities and
7 organizational structures, as well as credibility)
- 8 - provide an outlook on future options for including event reports in research which are given
9 by rapid technical and publishing developments (linked data, open access, semantic search
10 options)

11

12 *4) At present, all quality dimensions have the same weight. However, it could be argued that*
13 *some dimensions are more relevant than others. This point should be better investigated or at*
14 *least identified as a priority for future research.*

15 The referee raises a very important point. In fact, the weighting scheme has been the most
16 discussed point amongst the authors too. There are many possible points of view on a
17 document's quality. Based on the framework by (Wang and Strong, 1996) the 4 quality
18 criteria (QC) are the main pillars that define the overall quality of information (in their case
19 data) from a users point of view. So, one option would be to give even weight to each of the
20 QC. This however also means that the scores reached in each of the dimensions per QC will
21 be averaged. The main argument for using the same weight for each of the dimensions was
22 that it is most reflective of the chosen task at hand and therefore the user's perspective of our
23 particular study. We accompany this choice by the notification that "It is important to note
24 that P is not meant to label a document as per se bad or good and any new task at hand will
25 yield its own quality results. It provides a measure to assess the overall quality of a report and
26 assists in creating an overview of the quality present in the material. At any instance, this
27 overall score needs to be accompanied by an analysis of scores reached in the single
28 dimensions or combinations of dimension in order to identify the contextual scope of the
29 document and its strengths and limitations. (pg. 153, lines 9-14)".

1 The way to proceed and therefore a field for further research will be a user survey in order to
2 define those quality dimensions/categories that are most relevant and in order to derive any
3 weights.

4 We will stress this important aspect more clearly and add it to the concluding chapter as a
5 field of further research.

6

7 **1.3 Specific comments**

8 **1.3.1 Abstract**

9 *Comment 1: page 144 lines 23-25*

10 *“Using an example flood event that occurred in October/ November 1998 we demonstrate*
11 *how the information from multiple reports can be synthesised under consideration of their*
12 *quality”. This is not done in the paper. In section 3.4 there is not any consideration of the*
13 *quality of reports and their role in the overall information credibility. It’s just one main*
14 *criticism highlighted in previous section.*

15 We adapt this sentence according to the answer provided for major criticism no.2. I.e. we
16 rephrase to: “Using an example flood event that occurred in October/ November 1998 we
17 demonstrate the information from multiple reports can be synthesised.”

18 **1.3.2 Introduction**

19 *Comment 1: page 145 lines 7*

20 *What do you mean with “any systemic approach”? Any systemic approach to what?*

21 We rephrase: “any systematic event analysis”

22 *Comment 2: page 145 lines 15*

23 *“Contextual depth” is extensively defined in the following but, at this point of the paper, its*
24 *meaning is not clear to a wide audience. Please specify.*

25 We rephrase this sentence avoiding the usage of the term “contextual depth” in order to avoid
26 confusion with the term that is later on used in a defined way. Instead of “However, the
27 expectation towards the contextual depth of the documents seems to be rather unclear” we

1 rephrase: “However, the type and detail of information contained in the documents seems to
2 be rather unclear”

3 *Comment 3: page 147 lines 17-7*

4 *“They are commonly applied in the course of systematic reviews and meta-analyses are used*
5 *to synthesize the available evidence for a given question to identify and assess consistent*
6 *findings across diverse studies (i.e. statistical analysis of causal linkages, effectiveness of*
7 *interventions) and to inform policy (Burton, 2010; Borenstein et al., 2009)”. Not clear, please*
8 *rephrase.*

9 We will rephrase the entire paragraph from page 146, line 27 to page 147, line 7

10 “Evidence-based evaluations aim at synthesizing the available evidence for a given question
11 (e.g., how effective are interventions in a river system for habitat restoration of species x) to
12 identify and assess consistent findings across diverse studies and to inform policy (Burton,
13 2010; Borenstein et al., 2009). They are most commonly applied in the course of systematic
14 reviews and meta-analyses and have become standard in the health and medical sciences
15 (Higgins and Green, 2011) and have also been transferred to environmental science and
16 management (Centre for Evidence-Based Conservation, 2010; Norris et al., 2008; Osenberg et
17 al., 1999). Beside the quantitative meta-analyses that provide a reproducible weighted average
18 of the estimate of an effect, qualitative criteria-based methods of causal inference have been
19 developed (see Weed, 2000 for a comparison of both methods).”

20

21 *Comment 4: page 147 lines 23-24*

22 *What is “the environmental level (depicting the general Zeitgeist)”? Not clear*

23 This is a term used in historic hydrology/climatology, however, its use is probably not so
24 wide spread. We paraphrase this term rewriting the bracket to: “(depicting the general way of
25 thinking and expression during an epoch)”

26 *Comment 5: page 147 lines 29*

27 *As for contextual depth, “intrinsic quality assessment” is explained later in the paper and its*
28 *meaning is not so clear at this point. Please, specify*

29 We remove the word ‘intrinsic’ from this sentence as we refer to the general quality and
30 therefore need no specification here.

1 1.3.3 Methodology

2 *Comment 6: page 149 lines 22-23*

3 *“The spatial, temporal and contextual frame for the search is given by the task above”. Not*
4 *clear, specify.*

5 The sentence is repetitive to what is explained afterwards. We therefore is delete it.

6 *Comment 7: page 151 lines 26-28*

7 *“Within each of these dimensions the original contextual dimensions of Wang and Strong*
8 *(1996) are inherently considered”. Not clear, please specify.*

9 We include the dimensions in brackets so that the sentence reads: “Within each of these
10 dimensions the original contextual dimensions (added value, relevancy, completeness,
11 appropriate amount of information) of Wang and Strong (1996) are inherently considered”

12 *Comment 8: page 153 lines 4-6*

13 *“Assuming an average score QD_i of 0, 1, 2, or 3 over all dimensions (example: an average*
14 *score of 2 would 5 result in a score sum of $10 \times 2 = 20$ and $P = 20/30 = 0.67$), P can be*
15 *interpreted according to the quality labels of no, low, medium and high quality”. How ranges*
16 *for quality labels have been defined is not clear from this explanation. Please clarify*

17 We rephrase: The measure P can be interpreted in terms of quality labels, i.e. a document
18 being of no, low, medium and high quality. The ranges of P are based on the consideration of
19 an average score in all quality dimensions QD_i and the breaks are defined by the sum of
20 scores reached by the QD .

21 *Comment 9: page 153 lines 22-23*

22 *“In defining the quality dimensions we consider the spatial scope at which the report*
23 *documents an event as reference for the quality expectation and assessment”. How the spatial*
24 *scale plays on reports quality is not clear, even in the following application (section 3.4).*
25 *Please clarify*

26 In most cases the extent of a flood event will be larger than the one that is described in an
27 event report (national and federal states or local restrictions). At this point in the paper we
28 want to highlight, that we will assess the quality of each report with respect to its own spatial
29 scope and not with respect to whether the reports is reflective of the entire flood extent.

1 *Comment 10: page 153 lines 23-25*

2 *“In the document specific categories we introduce the convention that each report is assessed*
3 *with respect to its objectives”. Not clear, please specify*

4 Similar to the previous comment. Each report has its own objectives, i.e. one report maybe
5 only focussing on the meteorological aspects, another only on damages, and another on the
6 entire event. We want to evaluate the quality of each report in the quality that it reaches with
7 respect to its own specific objectives.

8 **1.3.4 Results**

9 *Comment 11: page 162 lines 1-2*

10 *“In the following we will discuss the agreements reached in the single QD with respect to the*
11 *strengths and weaknesses of the definitions of the dimensions and their respective grades”.*
12 *Not clear, please specify*

13 We reduce the sentence to “In the following we will discuss the agreements reached in the
14 single QD”, as this is the simple task we are about to perform.

15 *Comment 12: page 162 lines 12-14*

16 *“Differences in the assignments can be attributed to the large amount of variables that are*
17 *covered in any of the contextual QD which introduces a minor degree of subjectivity of a peer*
18 *in drawing the distinctions”. Why should more variables reduce subjectivity? Not clear,*
19 *please specify:*

20 We mean the large amount of technical aspects that are covered within each dimension of the
21 contextual quality category. We rephrase this sentence to: “Differences in the assignments are
22 the result of some degree of subjectivity in drawing the distinctions between the score classes
23 which will depend on the technical experience of a peer in any of the specific aspects like e.g.
24 hydro-meteorology.”

25 *Comment 13: page 164 lines 3-6*

26 *“In order to assess the effect of peer disagreement on the overall pedigree we compare the*
27 *resulting P values (Fig. 2). The maximum difference encountered is +0.13 equalling a score*
28 *difference of four (a difference of one score leads to an alteration of P by 0.03 units).” The*

1 *meaning of Fig. 2 is not clear. Relation between P difference and score difference is not*
2 *clear. Please, specify.*

3 Per document that was included in the peer review process, Fig. 2 plots the Pedigree scores
4 that were given by the author against the pedigree scores that were given by the Peers. P is the
5 ratio of the total sum of scores given per quality dimension divided by the maximum sum of
6 scores possible (in our case: 10 dimension * max 3 = 30). The plot and example highlight that
7 the peers and authors result for the quality assessment of any document are very close.

8 *Comment 14:*

9 *ISI journals cannot be considered grey literature. The proposed QAF can be used both to*
10 *evaluate grey and official literature. That is fine but must be clarify earlier in the paper.*

11 We will include this as a notion in the methodology section.

12 *Comment 15:*

13 *Most of discussed results are not evident form Table 3 or Figure 4 but supplementary*
14 *material is required. This should be highlighted.*

15 We add a note on that at the beginning of section 3.2.

16 *Comment 16 page 168 lines 15-16:*

17 *“Figure 4 shows (...) a pair wise correlation with the score class 3 of the contextual*
18 *dimensions and accuracy”. Not clear, please clarify.*

19 We rephrase to: “Figure 4 shows a clear correlation of both dimensions with the overall
20 quality of the documents. Those reports that are of an overall good quality are exclusively
21 well written and well structured.”

22 *Comment 17 page 168 line 20:*

23 *“83.5%”. Is it correct? According to the table the right value is 84.2%*

24 You are right. The percentage should be 84.2%.

25 *Comment 18 page 170 line 13:*

26 *“GDR”. What does it mean? Not defined before*

27 We define the abbreviation. GDR – German Democratic Republic.

28 *Comment 19 page 172 line 5:*

1 *“See section 3.2”. Reference is not correct.*

2 Reference to any section not needed here. Will be removed.

3 *Comment 20 page 174 line 3:*

4 *“ $Q(T < 5a)$ ”. Is it an error?*

5 It is correct. However we rephrase to: The main rivers were affected at a increasing gradient
6 south-north, with the upper and middle Rhine experiencing peak flow of small return periods
7 $Q(T < 5a)$ (#148, #28) and higher peak flows with increasing contributions from tributaries
8 Neckar, Main, Moselle.

9 **1.3.5 Discussion**

10 *Comment 21 page 176 lines 14-16:*

11 *“A natural extension of the example application presented is the combination of data based*
12 *and model-based results with the quality-labelled information of the reports resulting*
13 *essentially in an uncertainty assessment of the available knowledge”. This seems a very*
14 *important point but is not clear. Please, rephrase and clarify.*

15 *Comment 22 page 176 lines 20-22:*

16 *“Evidence-based or related methods are a natural successor of the results of this study that*
17 *can assist in combining quantitative and qualitative measures of uncertainty”. This seems a*
18 *very important point but is not clear. Please, rephrase and clarify.*

19 In the paragraph related to by comments 21 and 22 we want to highlight that our study
20 provides a starting point for an improved understanding of flood events. In our case we use
21 reports and provide a quality assessment scheme. Further research will be required to develop
22 a framework to combine these sources of information with results from model or data based
23 analysis. Possible frameworks can be the information expansion scheme provided by (Merz
24 and Blöschl, 2008) or evidence based methods like that of (van der Sluijs et al., 2005) or
25 (Norris et al., 2008).

26 **1.4 Technical corrections**

27 *Page 159 line 7: “bijective”*

28 *Page 167 line 2: Fig. 4 is the right one*

1 *Table A1: “efinitions”*

2 The errors will be corrected.

3

4 **References**

5 Merz, R., and Blöschl, G.: Flood frequency hydrology: 2. Combining data evidence, *Water*
6 *Resources Research*, 44, W08433, doi:08410.01029/02007WR006745, 2008.

7 Norris, R., Nichols, S. J., Ransom, G., Webb, A., Stewardson, M., Liston, P., and Mugodo, J.:
8 Causal criteria analysis. *Methods manual: a systematic approach to evaluate causality in*
9 *environmental science*, eWater Cooperative Research Centre, Canberra, 2008.

10 van der Sluijs, J. P., Craye, M., Funtowicz, S., Kloprogge, P., and Ravetz, J.: Combining
11 quantitative and qualitative measures of uncertainty in model-based environmental
12 assessment: The NUSAP system, *Risk analysis*, 25, 481-492, 2005.

13 Wang, R. Y., and Strong, D. M.: Beyond accuracy: what data quality means to data
14 consumers, *Journal of Management Information Systems*, 12, 5-33, 1996.

15

16