

We thank the reviewer for the time spent in reviewing our paper and making very helpful suggestions. We provided a point-by-point response to the reviewer's comments.

Anonymous Referee #1

Overall I think the paper is very interesting and generally well presented. The topic covered is quite complex (many different components to the modeling study) and therefore it is important that the explanations are as clear as possible. In most cases, I think this is true...

Author's response

We would like to thank the reviewer for his interest in this work.

Anonymous Referee #1

... but there is one example where I think the explanation could be improved and that is in the last paragraph of 4.1.2 where the three models are discussed. I think it needs to be made explicit that models 6, 30 and 54 are linked to the three different snow accumulation schemes.

Author's response

Agreed.

Author's changes in the manuscript

The following statement on page 12159, lines 6–7 of the discussion paper:

“... obtained with three competing model hypotheses (no. 6, 30 and 54) differing only in their snowmelt-accounting options.”

Has been replaced in the updated manuscript with the following statement:

“... obtained with three competing model hypotheses (no. 6, 30 and 54) differing only in their snowmelt-accounting options (respectively B1a, B1b and B1c).”

Anonymous Referee #1

I also found the explanation at the start of 4.2.3 to be quite confusing and could be explained a little better.

Author's response

Agreed. We hope that the following changes in the updated manuscript will help clarify our statement.

Author's changes in the manuscript

The following statement on page 12160, lines 24–26 of the discussion paper:

“The high representation of options F2a and F2b in Cluster 1 suggests that the catchment actually behaves as a serial system and may reveal a better correspondence with its overall physical structure.”

Has been replaced in the updated manuscript with the following statement:

“The frequency of options F2a and F2b in the best-performing cluster suggests that the catchment actually behaves as a ‘serial’ system.”

Please note that this sentence has been removed from Section 4.2 to Section 5.

Anonymous Referee #1

I found the implied definition of equifinality on page 12163 to be very limited. Why is equifinality limited to a single criterion? The concept was borrowed from geomorphology and relates to the same outcome from different causative processes. The definition used in the paper is a very limited ‘mathematical/statistical’ one.

Author's response

As underlined by the referee, the concept of equifinality relates to the same outcome from different causative processes. On Page 12163 line 2–4, we wrote that “two parameter sets are said to be equifinal if they can be regarded as equally acceptable in a statistical sense with respect to one particular criterion”. It seems to us that if one replaces the words “criterion” and “parameter sets” in our sentence by, respectively, “outcome” and “different causative processes”, one gets the original meaning of the concept given by the referee. We slightly modified our sentence to make more explicit that the concept of equifinality is defined here in a statistical context and not in general terms.

Author's changes in the manuscript:

The following statement on page 12159, lines 6–7 of the discussion paper:

“...two parameter sets are said to be equifinal if they can be regarded as equally acceptable in a statistical sense with respect to one particular criterion.”

Has been replaced in the updated manuscript with the following statement:

“...two parameter sets are said to be equifinal in a statistical sense if they can be regarded as equally acceptable with respect to a given model outcome.”

Anonymous Referee #1

It might have been useful to show some time series of flow and rain at the start of the paper to illustrate the hydrological regime (2.3.2). This could help the readers to understand the concepts of greater than 100% runoff coefficients. I assume that these are related to quite slow groundwater release processes where precipitation (or snowmelt) from one year only appears as runoff in the following year. Perhaps this also depends on how you define the hydrological year and this is not adequately explained in the paper.

Author's response

We thank the referee for this relevant suggestion. The hydrological year was defined from May to April so as to capture the snowmelt and peak flow seasons at mid-year. As explained in Sect. 2.3.2, these values of runoff coefficients were most likely due to an underestimation of precipitation at high elevations or to “a greater contribution of groundwater to surface flow”. We realized that this statement was not clear enough and modified it as indicated below.

Author's changes in manuscript

A new figure representing multi-decadal hyetograph and hydrograph has been added to the manuscript to illustrate the hydrological regime of the catchment studied. The definition of the hydrological year was inserted in the caption of this figure (Figure 2 in the updated manuscript). Moreover, the following statement on page 12146, line 22 of the discussion paper:

“... or a greater contribution of groundwater to surface flow...”

Has been replaced in the updated manuscript with the following statements:

“... or a delayed contribution of groundwater to surface flow from one year to another...”

Anonymous Referee #1

While the authors introduce some ‘real hydrology’ in section 4.2 these discussions are quite limited compared to the much greater detail about the statistics and mathematics of uncertainty. This aspect of the paper could be improved.

Author's response

We do agree that Section 4.2 is quite limited compared to the other parts of the discussion paper. This is mainly because we wish to limit these statements to very basic assumptions requiring much caution, given the lumped conceptual nature of the models involved.

Author's changes in the manuscript

Please note that these comments have been removed from Section 4.2 and put in Section 5 to emphasize their hypothetical nature.

Anonymous Referee #1

I also noted that the issues of data uncertainty associated with the estimation of natural streamflow are only mentioned right at the end, while these could have a very large impact on the modelling results if the naturalization process and the knowledge of abstractions is poor.

Author's response

Agreed. This comment was also made by the other anonymous referee. We admit that this point was not made clear in the paper and this was mainly due to space limitations. As explained in Section 2.1., vineyards and orchards cover most of the valley floors and lower hill slopes, where they benefit from a unique combination of clear skies, high temperatures and overall dry conditions throughout the growing season. Most of the annual precipitation, however, occurs as snow during the winter months, leading to an entire dependence on surface-water resources to satisfy crop water needs during the summer. Irrigation water abstractions occur at multiple locations along the river's course depending on both historical water rights and water availability. Because these abstractions are likely to influence the hydrological behavior of the catchment, especially during low-flow periods, they were added back to the observed streamflow before calibrating the models. This inevitably adds some uncertainty to the modeling of daily stream flows because a significant part of surface-water abstractions actually return to the river system within a few days. In general, ignoring these return flows will lead to overestimating natural stream flows on a daily basis. In this paper, however, the actual water withdrawals were not known with precision but only as percentages of the nominal water rights (these percentages are fixed on a monthly basis by the authorities depending on water availability), so the overall effects of streamflow naturalization on model uncertainty remained unknown.

Author's changes in manuscript

The following statements on page 12143, lines 17–18 of the discussion paper:

“...but account for less than 1% of the total catchment area (INE, 2009; CIREN, 2011). By contrast, natural vegetation outside the valleys is extremely sparse...”

Has been replaced in the updated manuscript with the following statements:

“...but account for less than 1% of the total catchment area (INE, 2009; CIREN, 2011). Most of the annual precipitation, however, occurs as snow during the winter months, leading to an entire dependence on surface-water resources to satisfy crop water needs during the summer. Irrigation water abstractions occur at multiple locations along the river's course depending on both historical water rights and water availability. By contrast, natural vegetation outside the valleys is extremely sparse...”

The following statement on page 12144, lines 22–25 of the discussion paper:

“Naturalized streamflow time series were estimated using information provided by the Chilean *Dirección General de Aguas*, mainly streamflow measurements at the gauging station of Rivadavia and historical surface-water diversion data.”

Has been replaced in the updated manuscript with the following statements:

“Water abstractions for irrigation were estimated using information on historical water allocations provided by the Chilean authorities. Because these abstractions are likely to influence the hydrological

behavior of the catchment during recession and low-flow periods, they were added back to the gauged streamflow in Rivadavia before calibrating the models.”

The following statements on page 12164, lines 21–27 of the discussion paper:

“It was also possible to highlight some errors in the streamflow data. Part of these errors might be associated with uncertainties in the estimation of natural streamflow. Further research is therefore required to better integrate the effect of water abstractions in the hydrological modeling process. From a multiple-hypothesis perspective, the modeling of irrigation water withdrawals should be regarded as a testable model component in its own right.”

Have been replaced in the updated manuscript with the following statements:

“It was also possible to highlight some errors in the streamflow data. The observed streamflow was ‘naturalized’ by simply adding back the estimated historical water abstractions. When applied on a daily basis, this process inevitably adds some uncertainty because a significant part of surface-water abstractions actually return to the river system within a few days due to conveyance and field losses. In general, ignoring these return flows leads to overestimating daily natural flows. In this paper, however, the actual water withdrawals were not known with precision but only as percentages of the nominal water rights (these percentages being fixed on a monthly basis by the authorities to account for water availability), so the overall impact of streamflow naturalization on model uncertainty remained unknown. Further research is underway to integrate the effect of water abstractions and crop water-use in the hydrological modeling process (Hublart et al., 2015; see also Kiptala et al., 2014). From a multiple-hypothesis perspective, the modeling of irrigation water water-use should be regarded as a testable model component in its own right.”

Anonymous Referee #1

I think the paper contains too many references - it is not a review paper and many of them are somewhat superfluous. There are also several that are included in the reference list that are not used in the text (Clark et al, 2009; Fenicia et al., 2007; Fowler and Kilsby, 2007; Freer et al., 2013; Hrachowitz et al., 2013; Krueger et al., 2010; Lang and Braun, 1990; Leavellesley et al., 2002; Loukas et al., 2002; Montecinos and Patricio, 2003; Olssen and Andersson, 2007; Staudinger et al., 2011; Strauch et al, 2006 and Zhang et al., 2010). Some of these could be related to wrong dates as the following included in the text could not be found in the list: Clark et al., 2005; Fenicia et al, 2006; Freer et al, 2003; Montecinos and Aceituno, 2003). Shaefli et al, 2011 is also spelt wrong and Souvignet et al. has the wrong date?

Author's response

We apologize for all these typos which were corrected in the updated manuscript. Moreover, only the most relevant references have been kept in the revised paper..

Author's changes in manuscript

Agreed. We apologize for these typos which have been corrected in the updated manuscript.

Anonymous Referee #1

Figures 4 to 8 could all be improved in clarity with larger font sizes and other improvements. There is space to do this.

Author's response

Agreed.

Author's changes in manuscript

These figures have been modified with larger front sizes and minor modifications in the updated manuscript.

Some minor points:

Anonymous Referee #1

Are the 12 and 8 (precip & temp) stations supposed to be shown on Figure 1?

Author's response

No, the weather stations could not be shown on Figure 1 because many of them are actually located outside the catchment.

Author's changes in manuscript

In the updated manuscript, Figure 1 has been modified to include those precipitation and temperature stations which belong to the catchment.

Anonymous Referee #1

Page 12149 line 14 – where is Eq 1 referred to?

Author's response

We apologize for this typo. This reference to “Eq . 1” is completely undue and was removed from the paper in the updated manuscript.

Anonymous Referee #1

Page 12157 line 8 – Is ‘emblematic’ the right word here’?

Author's response

We used “emblematic” in the sense of “illustrative” or “representative” but, as non-native English speakers, we cannot be totally sure of this choice. In the updated manuscript, this adjective was simply removed without impacting the overall meaning of the sentence.

Anonymous Referee #1

Page 12159 line 6 – ‘... internal state variable obtained...’

Author's response

Agreed and modified.

Anonymous Referee #1

Page 12160 line 5 – ‘... absence of sublimation...’

Author's response

Agreed and modified.

Anonymous Referee #1

Page 12161 – The reference of Figure 7 at the start of 4.3 should be Figure 8 I presume.

Author's response

Agreed and modified.

Anonymous Referee #1

Page 12162 line 14 – ‘...filling of a moisture...’

Author's response:

Agreed. Please note that this sentence is no longer used in the updated manuscript.

We thank the reviewer for the time spent in reviewing our paper and making very helpful suggestions. We provided a point-by-point response to the reviewer's comments.

Anonymous Referee #2

The article addresses the interesting issue of structural uncertainty in conceptual hydrological modelling. The authors test a large number of alternative structures on a catchment in the Andes in Chile and discuss their relative merits in a multi-objective framework.

Overall, I found the article interesting and well written. I think it could make a valuable contribution to HESS provided that a number of points are improved. I have two main concerns. First, the conclusions of this study do not appear so novel compared to existing works based either on multi-hypotheses or multi-objective frameworks. I think the authors should strengthen the last part (discussion/conclusion) of their paper to better demonstrate what was learnt from the quite complex testing scheme they set up and what is new compared to what was already shown in past studies.

Author's response

We thank the reviewer for this important remark. Our paper effectively draws upon the combination of a modular multiple-hypothesis approach with a multi-objective optimization scheme. We did not claim that these two modeling strategies were new (although modular approaches to multiple-hypothesis testing remain rare in comparison with multi-model approaches), but that the potential benefits of combining them within the same framework remain largely unexplored in current studies.

Perhaps more importantly, our study addresses a current lack of hydrological modeling effort in semi-arid Andes. As mentioned in the introduction part of the paper, "very few catchments in this region have been studied intensively enough to provide reliable model simulations, often with no estimation of the surrounding uncertainty". These two points are further detailed below in our answers to the following reviewer's comments. They have been emphasized in the new version of the manuscript.

Author's changes in manuscript

The discussion part of our paper has also been rewritten to better demonstrate what was learnt from the testing of a large number of model structures. In particular, note that the hypotheses made in Section 4.2 regarding possible links between model structures and the physical features of the catchment have been transferred to Section 5 ("Discussion and conclusion") in the updated manuscript.

Anonymous Referee #2

Second, their study would give more general conclusions if tests had been made on more than one catchment. Indeed, the conclusions may strongly depend on the characteristics of the selected catchment. It would be useful to test the approach to at least another catchment, to check whether similar conclusions are reached.

Author's response

In general, we agree that multiple-hypothesis frameworks should be tested on several catchments if one wishes to identify possible links between the physical features of a given catchment and some specific modeling decisions. This would also be very desirable considering the influence of data errors on the results obtained with any particular model structure or performance measure. While the need for comparative studies was only briefly mentioned in the paper (p. 12164, lines 9–11), it is further discussed in the updated manuscript (see Section 5, "Discussion and conclusion").

However, we would like to emphasize the fact that this study represents the first step of a larger research project, whose final aim is to assess the capacity to meet current and future irrigation water requirements in the Claro River catchment. Because considerable time and effort had to be devoted to gathering/interpolating the input data and implementing/testing the modeling framework, it was also necessary to limit the scope of our study to this particular catchment. Moreover, the main objective here was not to establish unambiguous relationships between the physical characteristics of Andean catchments and specific model requirements, but rather to assess the uncertainty associated with model non-uniqueness and structural inadequacy in the Claro River catchment. From this point of view, it should be stressed that the paper already provides a reliable framework by testing a total of 72 competing model structures in a region where catchment-scale conceptual models remain largely under-used. Adding other Andean catchments would be of particular interest to the objectives of the study if precipitation data on these catchments were available over the same 30-year period and could be considered more reliable. To our knowledge, this is not

the case in the Andes in general. The dataset used in the paper actually includes several of the highest weather stations available at this time scale in the Chilean Andes.

Author's changes in manuscript

We added a few comments to qualify our statements and insist on the need for comparative studies to confirm the hypotheses made in Section 4.2 regarding possible links between model structures and the physical features of the catchment. As mentioned above, note that these hypotheses have also been transferred to Section 5 (“Discussion and conclusion”) in the updated manuscript.

Anonymous Referee #2

I have also a number of detailed comments below. I think the paper could be reconsidered for publication after major revision.

Detailed comments:

1. Anonymous Referee #2

There are remaining typos that should be corrected. Consistency between references in the text and the list of references at the end of the manuscript should also be further checked.

Author's response

Agreed. We apologize for these typos which have been corrected in the updated manuscript.

Author's changes in manuscript

Please see the updated manuscript for the detailed corrections.

2. Anonymous Referee #2

Page 12139, line 25: The authors may find interesting reflections on this issue in the book edited by Wainwright and Mulligan (2004).

Author's response

We thank the referee for this interesting suggestion which has been inserted in the updated manuscript.

Author's changes in manuscript

The following statement on page 12139, lines 27–28 of the discussion paper:

“... as ready-made engineering tools with little or no consideration for the specific features of each catchment (Savenije, 2009)”

Has been replaced in the updated manuscript with the following statement:

“... as ready-made engineering tools with little or no consideration for the specific features of each catchment (Wainwright and Mulligan, 2004; Savenije, 2009)”

3. Anonymous Referee #2

Page 12142, lines 1–10: I do not agree that the multi-model approach was mainly focused on small catchments. There are a number of studies in the literature that investigated larger ranges of catchment size.

Author's response

Here we think there may be a definitional issue. In our opinion, a substantial distinction should be made between current multi-model strategies and modular modeling frameworks (MMF). While both rely to some extent on the concept of multiple-hypothesis testing, it should be noted that modular approaches offer the additional opportunity to examine the effect of each individual hypothesis (i.e. each modeling decision) by modifying only one component or constitutive equation at a time and testing a wide range of alternative combinations between model components. By

contrast, multi-model strategies generally involve ready-made model structures borrowed from the literature (“off-the-shelf” models).

We think this distinction was made clear on Page 12141 (lines 25–27) but maybe not enough on the following page where the issue of catchment size was discussed. Generally speaking, multiple-hypothesis frameworks should not be completely identified with modular modeling frameworks (as we did on Page 12142 and later in the paper) since in reality the latter represents only one possible approach to multiple-hypothesis testing. The manuscript was therefore modified to maintain a distinction between the two. This is important because, to our knowledge, most conceptual modular frameworks currently used in hydrological modeling studies have been applied to relatively “small” catchments of, at most, a few hundreds of km². As argued by the reviewer, this is not the case of more traditional multi-model approaches, which indeed cover a much larger range of catchment sizes.

To our knowledge, there is only one example of a study making use of a modular framework on a semi-arid catchment of more than 2000 km² and that is the original paper of Clark et al. (2008) introducing the FUSE toolbox. We agree that this point was somewhat overlooked in our paper and modified this part of the introduction to balance our statements with other arguments.

Author’s changes in the manuscript

The following statements on page 12143, lines 17–18 of the discussion paper:

“So far, however, this method has mostly been applied to small (<10 km²) experimental (well-monitored) catchments (e.g. Clark et al., 2008; Smith and Marshall, 2010; Buytaert and Beven, 2011; McMillan et al., 2012b; Fenicia et al., 2014), with less attention being given to larger scales of interest (100–400 km²) (e.g. Kavetski and Fenicia, 2011; Coxon et al., 2013) or long time periods. Therefore, the need remains to establish whether MHF can also be used to improve conceptual modeling on multi-decadal periods at operational scales of 1000 km² or more. The potential benefits of combining MHF with Pareto-based optimization schemes also remain largely unexplored in the current literature.”

Have been replaced with the following statements:

“So far, however, this method has mostly been applied to relatively small (<500 km²) and humid catchments of the Northern Hemisphere (Krueger et al., 2010; Smith and Marshall, 2010; Staudinger et al., 2011; Kavetski and Fenicia, 2011; McMillan et al., 2012b; Coxon et al., 2013), with less attention being given to larger scales of interest (>1000 km²) and semi-arid regions (Clark et al., 2008). Moreover, several of these studies have insisted on the need for multiple criteria related to different aspects of the system’s behavior in order to improve the usefulness of MMF. Yet, most of the time these additional criteria or signatures were not used to guide model development or constrain calibration but rather as posterior diagnostics in validation (see e.g. Kavetski and Fenicia, 2011). Thus, the potential benefits of using the concept of Pareto-efficiency to constrain model development and help differentiate between a large number of competing hypotheses remain largely unexplored in the current literature devoted to MMF. Also, very few studies have included alternative conceptual representations of snow processes in their modular frameworks (e.g. Smith and Marshall, 2010), even though snowmelt may have played a significant role in several cases (Clark et al., 2008; Staudinger et al., 2011).”

Anonymous Referee #2

Besides, what makes the application of such approaches to larger catchments essentially different given the lumped approach used? I found that the argument of scale to explain the novelty of the study not really convincing here.

Author's response

We do agree that the argument of scale may not be the most relevant of all to explain the novelty of our study. More fundamentally, we would like to insist on the fact that very little is currently known about the adequacy of commonly used conceptual models in semi-arid Andean catchments. Most modular multiple-hypothesis frameworks such as FUSE, SUPERFLEX and RRMT have been applied to humid or subhumid catchments characterized by a limited role of snowmelt. Moreover, as mentioned in the introduction, there is no strong evidence in our opinion that lumped conceptual models designed on small catchments remain adequate at larger spatial scales.

Author's changes in the manuscript

As mentioned above, the argument of scale has been further qualified with other arguments in the updated manuscript.

4. Anonymous Referee #2

Section 2: As explained above, I found that adding another case study (possibly under similar or different conditions) would make conclusions more general. Here the catchment is quite specific in the sense that there seems to be a huge uncertainty in precipitation estimates. Adding another catchment with better known precipitation would provide a comparative reference to balance the results presented here.

Author's response

Please see our answer on the first page of this document.

5. Anonymous Referee #2

Page 12143, line 26: The location of gauges could be shown in Fig. 1.

Author's response

Agreed. This comment was also made by the other anonymous referee. Note however that all the weather stations cannot be shown on Figure 1 because many of them are actually located outside the catchment.

Author's changes in the manuscript

Figure 1 was modified to include those precipitation and temperature stations which belong to the catchment.

6. Anonymous Referee #2

Page 12144, lines 17–22: I did not understand why the Oudin's PE formula was adjusted to the Penman-Monteith's one. Why not directly using the latter if it is found more adapted to the study site?

Author's response

We did not find the physically-based Penman-Monteith approach more adapted to the objective of our paper. Oudin et al. (2005) showed that this approach may actually be less advantageous than more empirical formulas when used in daily conceptual models. This is why we chose to use the Oudin's formula in this study. The reason why we chose to adapt its parameters K_1 and K_2 to the local conditions is because our study was part of a larger project which involved the assessment of current and future irrigation water requirements in the catchment. In this project, the Penman-Monteith equation appeared more suited to simulating crop water needs in the valleys, but, because it required meteorological data that were only available for the last three or four years (relative humidity, wind speed, solar irradiance), it was decided to rely on a modified version of Oudin's temperature-based formula, in which the values of K_1 and K_2 were determined by selecting those giving the best fit to the available Penman-Monteith estimates of PE. In fact, these modified values of K_1 and K_2 were very close to those found by Oudin et al. (2005) and a sensitivity analysis showed that such modifications had no impact on the performance of the hydrological models used in the present paper. As a consequence, we kept these modified values to remain consistent with the other part of the project.

Author's changes in manuscript

In order to simplify our statement and avoid any misunderstanding, we removed these details on the estimation of PE from the updated manuscript. Instead, the reader is referred to Hublart et al. (2014) for more details on the values of K_1 and K_2 .

7. Anonymous Referee #2

Page 12144, lines 22–25: This statement is a bit vague. Could the authors give more details on this and explain to which extent the naturalization process may introduce uncertainty in the evaluation of models?

Author's response

Agreed. This comment was also made by the other anonymous referee. We admit that this point was not made clear in the paper and this was mainly due to space limitations. As explained in Section 2.1., vineyards and orchards cover most of the valley floors and lower hill slopes, where they benefit from a unique combination of clear skies, high temperatures and overall dry conditions throughout the growing season. Most of the annual precipitation, however, occurs as snow during the winter months, leading to an entire dependence on surface-water resources to satisfy crop water needs during the summer. Irrigation water abstractions occur at multiple locations along the river's course depending on both historical water rights and water availability. Because these abstractions are likely to influence the hydrological behavior of the catchment, especially during low-flow periods, they were added back to the observed stream flows before calibrating the models. This inevitably adds some uncertainty to the modeling of daily stream flows because a significant part of surface-water abstractions actually return to the river system within a few days. In general, ignoring these return flows will lead to overestimating natural stream flows on a daily basis. In this paper, however, the actual water withdrawals were not known with precision but only as percentages of the nominal water rights (these percentages are fixed on a monthly basis by the authorities depending on water availability), so the overall effects of streamflow naturalization on model uncertainty remained unknown.

Author's changes in manuscript

The following statements on page 12143, lines 17–18 of the discussion paper:

“...but account for less than 1% of the total catchment area (INE, 2009; CIREN, 2011). By contrast, natural vegetation outside the valleys is extremely sparse...”

Has been replaced in the updated manuscript with the following statements:

“...but account for less than 1% of the total catchment area (INE, 2009; CIREN, 2011). Most of the annual precipitation, however, occurs as snow during the winter months, leading to an entire dependence on surface-water resources to satisfy crop water needs during the summer. Irrigation water abstractions occur at multiple locations along the river's course depending on both historical water rights and water availability. By contrast, natural vegetation outside the valleys is extremely sparse...”

The following statement on page 12144, lines 22–25 of the discussion paper:

“Naturalized streamflow time series were estimated using information provided by the Chilean *Dirección General de Aguas*, mainly streamflow measurements at the gauging station of Rivadavia and historical surface-water diversion data.”

Has been replaced in the updated manuscript with the following statements:

“Water abstractions for irrigation were estimated using information on historical water allocations provided by the Chilean authorities. Because these abstractions are likely to influence the hydrological behavior of the catchment during recession and low-flow periods, they were added back to the gauged streamflow in Rivadavia before calibrating the models.”

The following statements on page 12164, lines 21–27 of the discussion paper:

“It was also possible to highlight some errors in the streamflow data. Part of these errors might be associated with uncertainties in the estimation of natural streamflow. Further research is therefore required to better integrate the effect of water abstractions in the hydrological modeling process. From a multiple-hypothesis perspective, the modeling of irrigation water withdrawals should be regarded as a testable model component in its own right.”

Have been replaced in the updated manuscript with the following statements:

“It was also possible to highlight some errors in the streamflow data. The observed streamflow was ‘naturalized’ by simply adding back the estimated historical water abstractions (Sect. 2.2). When applied on a daily basis, this process inevitably adds some uncertainty to streamflow values because a

significant part of surface-water abstractions actually return to the river system within a few days due to conveyance and field losses. In general, ignoring these return flows would lead to overestimating daily natural flows. In this paper, however, the actual water withdrawals were not known with precision but only as percentages of the nominal water rights – these percentages being fixed on a monthly basis by the authorities to account for variations in water availability. The combined impact of streamflow and precipitation errors on the assessment of structural uncertainty thus remained unknown. Further research is currently underway to integrate the effects of water abstractions and crop water-use in the hydrological modeling process (Hublart et al., 2015; see also Kiptala et al., 2014 for another approach). From a multiple-hypothesis perspective, the modeling of irrigation water water-use should be regarded as a testable model component in its own right.”

8. Anonymous Referee #2

Page 12146, lines 10–15: Is not there any seasonality in these processes?

Author's response

The referee is correct to raise the question of seasonal variations in sublimation processes. However, snow cover in the catchment is only present during the winter months and it seemed reasonable, as a first approximation, to assume that sublimation rates remain constant during this period.

9. Anonymous Referee #2

Page 12146, line 22: Do the authors mean that the geological boundaries may be different from the topographic ones?

Author's response

No, that is not what we meant. We hope that the following changes in the updated manuscript will clarify our statement.

Author's change in manuscript

The following statement on page 12146, line 22 of the discussion paper:

“... or a greater contribution of groundwater to surface flow...”

Has been replaced in the updated manuscript with the following statements:

“... or a delayed contribution of groundwater to surface flow from one year to another...”

10. Anonymous Referee #2

Page 12151, line 25: Do the authors wish to refer to Section 2.3.1 instead?

Author's response

Agreed. We apologize for this typo which has been corrected in the updated manuscript.

11. Anonymous Referee #2

Page 12152, line 21: It is unclear how the SCA was modelled given the lumped approach followed here.

Author's response

There seems to be a misunderstanding about how SCA data were used in our study. What was modelled by the snow-accounting options is the total snow water equivalent (SWE) stored in the catchment. As explained on page 12152 of the paper, the SCA data were used “to quantify the error made in simulating the seasonal dynamics” of snow processes “in terms of snow presence or absence” at the catchment scale. The snow error criterion described in Figure 3 corresponds to the number of days when SCA observations and SWE simulations disagree as to whether snow is present in the catchment (no matter ‘how much’ snow is present). It relies on an indirect comparison between these two quantities.

Author's changes in manuscript

To make this statement clearer, the following sentence on page 12152, lines 22–23:

“... were used to evaluate the consistency of snow-accounting modeling options in terms of snow presence or absence in the basin”

Has been changed into:

“... were used to evaluate the consistency of snow-accounting modeling options in terms of snow presence or absence at the catchment scale”

12. Anonymous Referee #2

Page 12154, lines 4–12: I found this choice questionable. Uncertainty bounds should refer to actual nominal values. For example, if one seeks to build 90% confidence intervals, then one should expect that the uncertainty bands contain 90% of the observations, not the maximum of observations. Does it mean here that the authors wish to build 100% confidence intervals? If one wishes to use other confidence intervals, how the approach should be applied?

Author's response

Here it seems necessary to clarify some choices. This paper aimed at assessing model inadequacy and non-uniqueness using a combination of two *non-probabilistic* approaches: a modular modeling framework and a multi-objective optimization scheme. More precisely, we aimed at assessing to which extent structural uncertainty could be reduced by identifying which minimal set of best-performing models maximized the number of observations covered by the ensemble of Pareto-envelopes. This strategy relies on the expectation that a maximum of observations should lie within the overall envelope, provided that structural uncertainty is adequately represented by the ensemble of Pareto-envelopes and that additional sources of uncertainty are negligible. The success or failure of this objective merely indicates to which extent the aforementioned assumptions are correct. Where the objective is to reach X% of the observations (with $X < 100$), one can modify the fourth step of the algorithm detailed on page 12154–12155 by considering only a fraction $X/100$ of $N_{\text{obs}}(N_{\text{max}})$ and changing the equality sign on Page 12155, line 4 into a greater-than (or equality) sign.

However, given the non-probabilistic nature of this approach, we have some serious reservations as to whether the resulting simulation bounds can be interpreted in terms of “confidence intervals” or “confidence bands”. By comparison with probabilistic methods, multi-objective schemes based on the concept of Pareto-efficiency do not provide any estimate of the residual error variance. The envelopes derived from the sets of Pareto-optimal solutions quantify only the uncertainty arising from the trade-offs between competing criteria and do not have a predefined statistical meaning. This is of course a major drawback of non-probabilistic approaches to uncertainty. However, more probabilistic methods based on the statistical description of model residuals also have their disadvantages. In our opinion, what this comment actually reflects is a common issue in hydrological modeling regarding the definition and assessment of structural uncertainty in probabilistic or non-probabilistic terms. We admit that investigating this issue was far beyond the scope of our study.

Author's changes in manuscript

Because we agree that these assumptions and choices made in defining structural uncertainty bounds may be questionable, we provided more information on their limitations in the discussion part (Section 5) of the updated manuscript. The following statements were inserted:

“Eventually, the number of models used to represent structural uncertainty was reduced by searching for the minimal set of best-performing structures which maximized the number of observations covered by the ensemble of Pareto-envelopes. It is important to make clear that model inadequacy and non-uniqueness were evaluated here in non-probabilistic terms. In particular, the Pareto-envelopes derived for each model structure quantify only the uncertainty arising from the trade-offs between competing criteria and do not have a predefined statistical meaning (Engeland et al., 2006). Consequently, the overall simulation bounds shown in Figure 8 cannot be easily interpreted as ‘confidence bands’. Although discussing the adequacy of non-probabilistic approaches to structural uncertainty was far beyond the scope of this study, it is interesting to analyze the reasons why between 15 and 20% of the observations remained outside the overall simulated envelope in both calibration and validation. To a large extent, this lack of performance can be attributed either to uncertainties in

the precipitation and streamflow data that were overlooked in this study or to an insufficient coverage of the hypothesis and objective spaces.”

Also, we modified the following statements on page 12154, lines 4–12 of the discussion paper:

“The overall uncertainty envelope should be wide enough to include most of the observed discharge but not so wide that its representation of the various aspects of the hydrograph (rising limb, peak discharge, falling limb, baseflow) becomes meaningless. In general, one will seek to reduce as much as possible the width of the envelope while maximizing the number of observations enclosed within the bounds. In this study, priority was given to maintaining at its lowest value the number of outlying observations before searching for the best combination of models which minimized the envelope area.”

These statements have been replaced with the following ones:

“The overall uncertainty envelope should be wide enough to include a large proportion of the observed discharge but not so wide that its representation of the various aspects of the hydrograph (rising limb, peak discharge, falling limb, baseflow) becomes meaningless. In this study, priority was given to maintaining at its lowest value the number of outlying observations before searching for the best combination of models which minimized the envelope area.”

Anonymous Referee #2

I understand that the authors rightly distinguish reliability and sharpness as two expected qualities of the uncertainty estimates, but there are many criteria proposed in the literature to evaluate these qualities. Maybe the authors should use the commonly applied criteria to strengthen the evaluation of uncertainty bounds.

Author's response

We are aware that many other criteria exist in the literature, in particular regarding sharpness. However, choosing one of these over the others seemed quite arbitrary in the absence of preliminary sensitivity analyses, for which we lacked time.

13. Anonymous Referee #2

Page 12156, lines 14–16: It is a bit difficult to see at first glance the structural differences between these three models. The reader has to reconstruct the structures from table 4 and figure 2. Could the authors help the reader here by detailing these differences?

Author's response

We thank the reviewer for his remark, which allowed us to clarify this point in the updated manuscript.

Author's changes in manuscript

The following sentence on page 12156, lines 14–16 of the discussion paper:

“Models no. 22, 46 and 54, for instance, yield very similar values of the high-flow criterion (Crit1), despite huge differences in their modeling options.”

Have been replaced with the following sentence:

“Models no. 22 (A1–B1a–C3–D2–E1–F2b), 46 (A1–B1b–C3–D2–E1–F2b) and 54 (A1–B1c–C1–D3–E2–F1b), for instance, yield very similar values of the high-flow criterion (Crit1), despite some differences in their modeling options.”

14. Anonymous Referee #2

Page 12162, lines 5–6: Was this actually demonstrated here, given there remains similarly performing structures? Besides, I think the usefulness of multi-model frameworks was already demonstrated by past studies. So maybe this should be seen more like a confirmation of existing results.

Author's response

We agree that “demonstrated” sounds a bit excessive here. It has been removed from the updated manuscript. Regarding the usefulness of multi-model frameworks, please see our answer on Page 1 of this document as well as the modifications provided in the updated manuscript.

15. Anonymous Referee #2

Page 12162, lines 16–22: Can 9-parameter models be considered as parsimonious? The difference between 9 and 13 parameters is not so large, since many modellers may consider 9-parameter models already overparameterized. Maybe this discussion could further refer to past works discussing parsimony in conceptual modelling.

Author's response

We thank the reviewer for helping us to further clarify the important issue of model parsimony in a multi-objective context. Many authors rightly consider that a maximum of 5 to 6 free parameters should be accepted in calibration when using a single objective function. Efstratiadis and Koutsoyiannis (2010) extended this empirical rule to the case of multi-objective schemes by allowing “a ratio of about 1:5 to 1:6 between the number of criteria and the number of parameters to optimize”. For a multi-objective scheme based on four criteria, this would lead to consider 20 to 24-parameter models as still being parsimonious, which, of course, would seem highly unlikely to many modelers. This is because in most cases, as Efstratiadis and Koutsoyiannis (2010) also pointed out, the various criteria used are not independent of each other. In our case, for instance, the information added by the low-flow criterion does not appear so different from that already introduced by the high-flow criterion. By contrast, the snow error criterion really adds new information on some specific snow-accounting parameters. Thus, 9-parameter models should not be regarded as being ‘parsimonious’ in general but only with respect to the number and quality of the criteria used in calibration.

Author's changes in manuscript

These reflections were included in the updated manuscript to clarify our statement (see Section 5).

16. Anonymous Referee #2

Page 12164, line 1: Would groundwater data be actually helpful in the case of this catchment, given the large uncertainties in precipitation estimates?

Author's response

The reviewer is correct to question the usefulness of additional groundwater data in our case. Additional information on precipitation would be probably far more relevant to improve the reliability of model predictions. This point was made clearer in the updated manuscript.

Author's changes in manuscript

The following words on page 12164, line 1 of the discussion paper:

“e.g. groundwater levels”

Have been replaced with:

“e.g. observed snow heights, irrigation water-use”

17. Anonymous Referee #2

Table 1: I do not understand the first equation for snow, which seems larger than P. Maybe remind the option type in the table.

Author's response

We apologize for this typo which has been corrected in the updated manuscript.

Author's changes in manuscript

Please see the modifications made to Table 1 in the updated manuscript.

18. Anonymous Referee #2

Table 2: Where does the range for K_c come from? The ranges given for K_3 seem dependent on the option but are the same in the table.

Author's response

The range of values tested for this parameter stem from the following assumption:

$$AE = K_{veg}Area_{veg}PE = K_CPE$$

where K_{veg} is a coefficient which varies between 0 and 1, and $Area_{veg}$ is the fraction of land covered with vegetation, which we limited to a maximum of 0.5 given the extreme aridity of the Claro River catchment.

Author's changes in manuscript

Initially, this explanation was not included in the discussion paper for brevity's sake. In the updated manuscript, we inserted a brief explanation on this point in the caption of Table 2.

Anonymous Referee #2

The ranges given for K_3 seem dependent on the option but are the same in the table.

Author's response

We apologize for this typo which has been corrected in the updated manuscript.

Author's changes in manuscript

Please see the modifications made to Table 2 in the updated manuscript.

Reducing structural uncertainty in conceptual hydrological modeling in the semi-arid Andes

P. HUBLART^{1,4}, D. RUELLAND², A. DEZETTER³ & H. JOURDE¹

¹UM2, ²CNRS, ³IRD – UMR HydroSciences Montpellier, Place E. Bataillon, 34395 Montpellier Cedex 5, France

⁴Centro de Estudios Avanzados en Zonas Áridas (CEAZA), Raúl Bitrán s/n, La Serena, Chile

paul.hublart@um2.fr / denis.ruelland@um2.fr

Abstract The use of lumped, conceptual models in hydrological impact studies requires placing more emphasis on the uncertainty arising from deficiencies and/or ambiguities in the model structure. This study provides an opportunity to combine a multiple-hypothesis framework with a multi-criteria assessment scheme to reduce structural uncertainty in the conceptual modeling of a meso-scale Andean catchment (1515 km²) over a 30-year period (1982–2011). The modeling process was decomposed into six model-building decisions related to the following aspects of the system behavior: snow accumulation and melt, runoff generation, redistribution and delay of water fluxes, and natural storage effects. Each of these decisions was provided with a set of alternative modeling options, resulting in a total of 72 competing model structures. These structures were calibrated using the concept of Pareto optimality with three criteria pertaining to streamflow simulations and one to the seasonal dynamics of snow processes. The results were analyzed in the four-dimensional space of performance measures using a fuzzy *c*-means clustering technique and a differential split sample test, leading to identify 14 equally acceptable model hypotheses. A filtering approach was then applied to these best-performing structures in order to minimize the overall uncertainty envelope while maximizing the number of enclosed observations. This led to retain 8 model hypotheses as a representation of the minimum structural uncertainty that could be obtained with this modeling framework. Future work to better consider model predictive uncertainty should include a proper assessment of parameter equifinality and data errors, as well as the testing of new or refined hypotheses to allow for the use of additional auxiliary observations.

1. INTRODUCTION

Conceptual catchment models based on the combination of several schematic stores are popular tools in flood forecasting and water resources management (e.g. Jakeman and Letcher, 2003; Xu and Singh, 2004). The main rationale behind this success lies in the fact that relatively simple structures with low data and computer requirements generally outweigh the performance of far more complex physically-based models (e.g. Michaud and Sorooshian, 1994; Refsgaard and Knudsen, 1996; Kokkonen and Jakeman, 2001). Also, most water management decisions are made at operational scales having much more to do with catchment-scale administrative considerations than with our understanding of microscale–fine-scale processes. As a result, conceptual models are being increasingly used to evaluate the potential impacts of climate change on hydrological systems (e.g. Minville et al., 2008; Ruelland et al., 2012) and freshwater availability (e.g. Milano et al., 2013; Collet et al., 2013).

This modeling strategy, however, is regularly criticized for oversimplifying the physics of catchments and leading to unreliable simulations when conditions shift beyond the range of prior experience. Part of the problem comes from the fact that model structures are usually specified *a priori*, based on preconceived opinions about how systems work, which in general leads to an excessive dependence on the calibration process. More than a lack of physical background, this practice reveals a misunderstanding about *how* such models should be based on physics (Kirchner, 2006; Blöschl and Montanari, 2010). Hydrological systems are not structureless things composed of randomly distributed elements, but rather self-organizing systems characterized by the emergence of macroscale patterns and structures (Dooge, 1986; Sivapalan, 2006; Ehret et al., 2014). As such, the reductionist idea that catchments can be understood by merely aggregating (upscaling) fine-scale mechanistic laws is generally misleading (~~Anderson, 1972~~; Dooge, 1997; McDonnell et al., 2007). Self-organization at the catchment scale means that new hydrologic relationships with fewer degrees of freedom have to be envisioned (e.g. McMillan, 2012a). Yet, finding simplicity in complexity does not imply that simple models available in the literature can be used as ready-made engineering tools with little or no consideration for the specific features of each catchment (Wainwright and Mulligan, 2004; Savenije, 2009). As underlined by Kirchner (2006), it is important to ensure that the “right answers” are obtained for the “right reasons”. In the case of poorly-defined systems where physically-oriented interpretations can only be sought *a posteriori* to check for the model realism, this requires placing more emphasis on the uncertainty arising from deficiencies and/or ambiguities in the model structure than is currently done in most hydrological impact studies.

57 Structural uncertainty can be described in terms of *inadequacy* and *non-uniqueness*. Model
58 inadequacy arises from the many simplifying assumptions and epistemic errors made in the selection
59 of which processes to represent and how to represent them. It reflects the extent to which a given
60 model differs from the real system it is intended to represent. In practice, this results in the failure to
61 capture all relevant aspects of the system behavior within a single model structure or parameter set. A
62 common way of addressing this source of uncertainty is to adopt a top-down approach to model-
63 building (Jothityangkoon et al., 2001; Sivapalan et al., 2003), in which different models of increasing
64 complexity are tested to determine the adequate level of process representation. Where fluxes and state
65 variables are made explicit, alternative data sources (other than streamflow) such as groundwater
66 levels (Seibert, 2000; Seibert and McDonnell, 2002), tracer samples (Son and Sivapalan, 2007; Birkel
67 et al., 2010; Capell et al., 2012) or snow measurements (Clark et al., 2006; Parajka and Blöschl, 2008),
68 can also be used to improve the internal consistency of model structures. Additional criteria can then
69 be introduced in relation to these auxiliary data or to specific aspects of the hydrograph (driven vs.
70 nondriven components, rising limb, recession limbs...). In this perspective, multi-criteria evaluation
71 techniques based on the concept of Pareto-optimality provide an interesting way to both reduce and
72 quantify structural inadequacy (Gupta et al., 1998; Boyle et al., 2000; Efstratiadis and Koutsoyiannis,
73 2010). A parameter set is said to be Pareto-optimal if it cannot be improved upon without degrading at
74 least one of the objective criteria. In general, meaningful information on the origin of model
75 deficiencies can be derived from the mapping of Pareto-optimal solutions in the space of performance
76 measures (often called the Pareto front) and used to discriminate between several rival structures (Lee
77 et al., 2011). Further, the Pareto set of solutions obtained with a given model is commonly used to
78 generate simulation envelopes (hereafter called 'Pareto-envelopes' for brevity's sake) representing the
79 uncertainty associated with structural errors (i.e. model inadequacy).

80 Non-uniqueness refers to the existence of many different model structures (and parameter sets)
81 giving equally acceptable fits to the observed data. Structural inadequacy and the limited (and often
82 uncertain) information of the available data make it highly unlikely to identify a single, unambiguous
83 representation of how a system works. There may be, for instance, many different possible
84 representations of flow pathways yielding the same integral signal (e.g. streamflow) at the catchment
85 outlet (Schaeffli et al., 2011). Non-uniqueness in model identification has also been widely
86 described in terms of equifinality (Beven, 1993 and 2006) and may be viewed as a special case of a
87 more general epistemological issue known as the “underdetermination” problem. Over the past
88 decade, these considerations have encouraged a shift in focus toward more flexible modeling tools
89 based on the concept of multiple working hypotheses (Buytaert and Beven, 2011; Clark et al., 2011).
90 A number of modular frameworks have been proposed, in which model components (i.e. individual
91 hypotheses) can be assembled and connected in many ways to build a variety of alternative model
92 structures (i.e. overall hypotheses). Recent examples of such modular modeling frameworks (MMF)
93 include the Imperial College Rainfall-Runoff Modeling Toolbox (RRMT) (Wagener et al., 2002), the
94 Framework for Understanding Structural Errors (FUSE) (Clark et al., 2008) and the SUPERFLEX
95 modeling environment (Fenicia et al., 2011). Clark et al. (2011) suggested that multiple-hypothesis
96 frameworks (MHF) this approach to model identification represents a valuable alternative to “most
97 practical applications of the top-down approach”, which “seldom consider competing process
98 representations of equivalent complexity”. Compared to current multimodel strategies, these
99 frameworks MMF also provide the possibility to better scrutinize the effect of each individual
100 hypothesis (i.e. model component), provided that the model decomposition is sufficiently fine-grained.
101 Finally, Clark et al. (2011) argued that ensembles of competing model structures obtained from MMF
102 (both of equal and varying complexity) can also be generated–used to quantify the structural
103 uncertainty arising because of system non-identifiability (i.e. model non-uniqueness). So far, however,
104 this method has mostly been applied to relatively small (<500 km²) and humid catchments of the
105 Northern Hemisphere (Krueger et al., 2010; Smith and Marshall, 2010; Staudinger et al., 2011;
106 Kavetski and Fenicia, 2011; McMillan et al., 2012b; Coxon et al., 2013), with less attention being
107 given to larger scales of interest (>1000 km²) and semi-arid regions (e.g. Clark et al., 2008). Moreover,
108 several of these studies have insisted on the need for multiple criteria related to different aspects of the
109 system’s behavior in order to improve the usefulness of MMF. Yet, most of the time these additional
110 criteria or signatures were not used to guide model development or constrain calibration but rather as
111 posterior diagnostics in validation (see Kavetski and Fenicia, 2011). Thus, the potential benefits of

112 using the concept of Pareto-efficiency to constrain model development and help differentiate between
113 numerous competing hypotheses remain largely unexplored in the current literature devoted to MMF.
114 Also, very few studies have included alternative conceptual representations of snow processes in their
115 modular frameworks (e.g. Smith and Marshall, 2010), even though snowmelt may have played a
116 significant role in several cases (Clark et al., 2008; Staudinger et al., 2011). So far, however, this
117 method has mostly been applied to small (<10 km²) experimental (well-monitored) catchments (e.g.
118 Clark et al., 2008; Smith and Marshall, 2010; Buytaert and Beven, 2011; McMillan et al., 2012b;
119 Fenicia et al., 2014), with less attention being given to larger scales of interest (100–400 km²) (e.g.
120 Kavetski and Fenicia, 2011; Coxon et al., 2013) or long time periods. Therefore, the need remains to
121 establish whether MHF can also be used to improve conceptual modeling on multi-decadal periods at
122 operational scales of 1000 km² or more. The potential benefits of combining MHF with Pareto-based
123 optimization schemes also remain largely unexplored in the current literature.

124 Addressing these issues is of particular importance in the case of arid to semi-arid, ~~mountainous~~
125 Andean catchments such as those found ~~in north-central Andes (30°S)~~ around 30°S. The Norte Chico
126 region of Chile, in particular, has been identified as being highly vulnerable to climate change impacts
127 in a number of recent reports (IPCC, 2013) and studies (e.g. Souvignat et al., 2010; Young et al.,
128 2010). Yet, very few catchments in this region have been studied intensively enough to provide
129 reliable model simulations, often with no estimation of the surrounding uncertainty (Souvignat,
130 ~~2007~~ 2008; Ruelland et al., 2011; Vicuña et al., ~~2012~~ 2011; Hublart et al., 2013). This study is the first
131 step of a larger research project, whose final aim is to assess the capacity to meet current and future
132 irrigation water requirements in a mesoscale catchment of the Norte Chico region. The objective here
133 is to provide a set of reasonable model structures that can be used for the hydrological modeling of the
134 catchment. To achieve this goal, a MHF-MMF was developed and combined with a multi-criteria
135 optimization framework using streamflow and satellite-based snow cover data.

137 2. STUDY AREA

138 2.1. General site description

140 The Claro River Catchment (~~CRC~~) is a semi-arid, mountainous catchment located in the
141 northeastern part of the Coquimbo region, in north-central Chile (Fig. 1). It drains an area of
142 approximately 1515 km², characterized by high elevations ranging from 820 m a.s.l. at the basin outlet
143 (Rivadavia) to over 5500 m a.s.l. in the Andes Cordillera. The topography is dominated by a series of
144 generally north-trending, fault-bounded mountain blocks interspersed with a few steep-sided valleys.

145 The underlying bedrock consists almost entirely of granitic rocks ranging in age from
146 Pennsylvanian to Oligocene and locally weathered to saprolite. Above 3000 m a.m.s.l., repeated
147 glaciations and the continuous action of frost and thaw throughout the year have caused an intense
148 shattering of the exposed rocks (Caviedes and Paskoff, 1975), leaving a landscape of bare rock and
149 screes almost devoid of soil.

150 The valley-fill material consists of mostly unconsolidated Quaternary alluvial sediments mantled
151 by generally thin soils (< 1 m) of sandy to sandy-loam texture (~~CIREN, 2005~~). Vineyards and orchards
152 cover most of the valley floors and lower hill slopes but account for less than 1% of the total
153 catchment area (~~INE, 2009; CIREN, 2011~~). Most of the annual precipitation, however, occurs as snow
154 during the winter months, leading to an entire dependence on surface-water resources to satisfy crop
155 water needs during the summer. Irrigation water abstractions occur at multiple locations along the
156 river's course depending on both historical water rights and water availability. By contrast, natural
157 vegetation outside the valleys is extremely sparse and composed mainly of subshrubs (e.g. *Adesmia*
158 *echinus*) and cushion plants (e.g. *Laretia acaulis*, *Azorella compacta*) with very low transpiration rates
159 (Squeo et al., 1993). The Claro River originates from a number of small tributaries flowing either
160 permanently or seasonally in the mountains.

161 2.2. Hydro-climatic data

163 | In order to represent the hydro-climate climatic variability over-of the catchment, a 30-year period
 164 | (1982–2011) was chosen according to data availability and quality. Precipitation and temperature data
 165 | were interpolated based on respectively 12 and 8 stations (Fig. 1) using the inverse distance weighted
 166 | method on a 5km x 5km grid. Since very few measurements were available outside the river valleys,
 167 | elevation effects on precipitation and temperature distribution were considered using the SRTM digital
 168 | elevation model (Fig. 1). In a previous study, Ruelland et al. (2014) examined the sensitivity of the
 169 | GR4j hydrological model to different ways of interpolating climate forcing on this basin. Their results
 170 | showed that a dataset based on a constant lapse rate of 6.5°C/km for temperature and no elevation
 171 | effects for precipitation provided slightly better simulations of the discharge over the last 30 years.
 172 | However, since the current study also seeks to reproduce the seasonal dynamics of snow accumulation
 173 | and melt, it was decided to rely on a mean monthly orographic gradient estimated from the
 174 | precipitation observed series (Fig. 1). Potential evapotranspiration (PE) was computed using the
 175 | following formula proposed by Oudin et al. (2005):
 176 |

$$PE = \frac{R_e}{\lambda \rho} \times \frac{T + K_2}{K_1} \quad \text{if } T + K_2 > 0 \quad \text{else } PE = 0 \quad (1)$$

177 |
 178 | where PE is the rate of potential evapotranspiration (mm.d⁻¹), R_e is the extraterrestrial radiation (MJ.m⁻².d⁻¹), λ is the latent heat flux (2.45 MJ.kg⁻¹), ρ is the density of water (kg.m⁻³), and T is the mean daily
 179 | air temperature (°C) and K₁ and K₂ are fitted parameters (for more details on the values of K₁ and K₂,
 180 | see Hublart et al. (2014)). Oudin et al. (2005) determined the values of K₁ and K₂ by selecting those
 181 | that gave the best streamflow simulations when the formula was used to feed hydrological models. In
 182 | this study, the FAO Penman Monteith equation for a reference grass was used as a basis to tune K₁ and
 183 | K₂ at two different locations within the basin (Rivadavia, Pisco Elqui, Fig. 1) (for more details on the
 184 | results, see Hublart et al. (2014)). Water abstractions for irrigation were estimated using information
 185 | on historical water allocations provided by the Chilean authorities. Because these abstractions are
 186 | likely to influence the hydrological behavior of the catchment during recession and low-flow periods,
 187 | they were added back to the gauged streamflow in Rivadavia before calibrating the
 188 | models. Naturalized streamflow time series were estimated using information provided by the Chilean
 189 | Dirección General de Aguas, mainly streamflow measurements at the gauging station of Rivadavia
 190 | and historical surface water diversion data. In addition to streamflow data, remotely-sensed data from
 191 | the MODerate resolution Imaging Spectroradiometer (MODIS) sensor were used to estimate the
 192 | seasonal dynamics of snow accumulation and melt processes over a 9-year period (2003–2011). Daily
 193 | snow cover products retrieved from NASA's Terra (MOD10A1) and Aqua (MYD10A1) satellites were
 194 | combined into a single, composite 500-m resolution product to reduce the effect of swath gaps and
 195 | cloud obscuration. The remaining data voids were subsequently filled using a linear temporal
 196 | interpolation method.
 197 |

198 | 199 | 2.3. Hydrological functioning of the catchment

200 | 201 | 2.3.1. Precipitation variability

202 | Among the primary factors that control the hydrological functioning of the CRC catchment is the
 203 | high seasonality of precipitation patterns. Precipitation occurs mainly between June and August during
 204 | the winter months when the South Pacific High reaches its northernmost position. Most of the annual
 205 | precipitation falls as snow at high elevations, where it accumulates in seasonal snow packs that are
 206 | gradually released from October to April. The El Niño Southern Oscillation (ENSO) represents the
 207 | largest source of climate variability at the interannual timescale (e.g. Rutllant and Fuenzalida, 1991;
 208 | Montecinos and Aceituno, 2003). Anomalously wet (dry) years in the region are generally associated
 209 | with warm (cold) El Niño (La Niña) episodes and a simultaneous weakening (strengthening) of the
 210 | South Pacific High. It is worth noting, however, that some very wet years in the catchment can also
 211 | coincide with neutral to weak La Niña conditions, as in 1984, while several years of below-normal
 212 | precipitation may not exhibit clear La Niña characteristics (Verbist et al., 2010; Jourde et al., 2011).
 213 | These anomalies may be due to other modes of climate variability affecting the Pacific basin on longer

214 timescales. The Interdecadal Pacific Oscillation (IPO), in particular, has been shown to modulate the
215 influence of ENSO-related events according to cycles of between 15 and 30 years (Schulz et al., 2011;
216 Quintana and Aceituno, 2012). Recent shifts in the IPO phase occurred in 1977 and 1998 and may be
217 responsible for the highest frequency of humid years during the 1980s and the early 1990s when
218 compared to the late 1990s and the 2000s.

219 2.3.2. Catchment-scale water balance and dominant processes

220 Notwithstanding this significant climate variability, a rough estimate of the catchment water
221 balance can be given for the period 2003–2011 using the data presented in the previous subsection and
222 additional information available in the literature. Spatially averaged precipitation ranges from a low of
223 80 mm in 2010 to an estimated high of 190 mm in 2008. Evapotranspiration from non-cultivated areas
224 is sufficiently low to be reasonably neglected at the basin scale (Kalthoff et al., 2006). By contrast,
225 water losses from the cultivated portions of the basin are likely to be around 10 mm.yr⁻¹ (Hublart et al.,
226 2013, 2014). At high elevations, sublimation plays a much greater role than evapotranspiration. Mean
227 annual sublimation rates over two glaciers located in similar, neighbouring catchments have been
228 estimated to be about 1 mm.d⁻¹ (see e.g. MacDonell et al., 2013). Thus, a first estimate of the annual
229 water loss associated with snow sublimation can be made by multiplying, for each day of the period,
230 the proportion of the catchment covered with snow by an average rate of 1 mm.d⁻¹. This leads to a
231 mean annual loss of 70 mm between 2003 and 2011. Note that this value is of the same order of
232 magnitude as those obtained by Favier et al. (2009) using the Weather Research and Forecasting
233 regional-scale climate model. Mean annual discharge per unit area varies from a minimum of 20 mm
234 in 2010 to a maximum of 140 mm in 2003. Interestingly, runoff coefficients exceed 100% during
235 several years of the period (in 2003, 2006, 2007 and 2009), indicating either an underestimation of
236 precipitation at high elevations, as suggested by Favier et al. (2009), or a greater-delayed contribution
237 of groundwater to surface flow from one year to another (Jourde et al., 2011).

238 Groundwater movement in the catchment is mainly from the mountain blocks toward the valleys
239 and then northward along the riverbed. In the mountains, groundwater flow and storage are controlled
240 primarily by the presence of secondary permeability in the form of joints and fractures (Souvignet
241 Strauch et al., 2006). The unconfined valley-fill aquifers are replenished by mountain front recharge
242 along the valley margins and by infiltration through the channel bed along the losing river reaches
243 (Jourde et al., 2011). Their hydraulic conductivity and saturated thickness range from about 10 m.d⁻¹
244 and 40 m respectively in the upper part of the catchment to more than 30 m.d⁻¹ and 60 m respectively
245 at the outlet (CAZALAC, 2006), allowing a rapid transfer of water to the hydraulically connected
246 surface streams. Pourrier et al. (2014) studied flow processes and dynamics in the headwaters of the
247 neighbouring Turbio River catchment; yet very little remains currently known about the emergent
248 processes taking place at the catchment scale.

249

250 3. METHODS

251

252 3.1. Multiple-hypothesis modeling framework

253 In order to evaluate various numerical representations of the catchment functioning, a multiple-
254 hypothesis modeling framework inspired by previous studies in literature was developed. All the
255 models built within this framework are lumped hypotheses run at a daily time step. The modeling
256 process was decomposed into three modules and six model-building decisions. Each module deals
257 with a different aspect of the precipitation–runoff relationship through one or more decisions (Fig. 2):
258 snow accumulation (A) and melt (B), runoff generation (C), redistribution (D) and delay (E) of water
259 fluxes, and natural storage effects (F). Each of these decisions is provided with a set of alternative
260 modeling options, which are named by concatenating the following elements: first a capital letter from
261 A to F referring to the decision being addressed, then a number from 1 to 3 to distinguish between
262 several competing architectures and, finally, a lower case letter from *a* to *c* to indicate different
263 parameterizations of the same architecture. Model hypotheses are named by concatenating the names

264 of the six modeling options used to build them (see Table 4). The models designed within this
265 framework share the same overall structure (based on the same series of decisions) but differ in their
266 specific formulations within each decision.

267 The model-building decisions can be divided into two broad categories. The first pertains to the
268 production of fluxes from conceptual stores (decisions B, C and F). The second concerns the
269 allocation and transmission of these fluxes using the typical junction elements and lag functions
270 (decisions A, D and E) described by [Fenicia et al. \(2011\)](#). Junction elements can be defined as
271 “zero-state” model components used to combine several fluxes into a single one (option D2) or split a
272 single flux into two or more fluxes (options A1 and D3). Lag functions are used to reflect the travel
273 time (delay) required to convey water from one conceptual store to another or from one or more
274 conceptual stores to the basin outlet. They usually consist of convolution operators (option E2),
275 although conceptual stores may also do the trick. Modeling options in which water fluxes are left
276 unchanged are labelled as “No operation” options in Fig. 2. Water fluxes and state variables are named
277 using generic names (from Q1 to Q6 and from S1 to S4, respectively) to ensure a perfect modularity of
278 the framework. Further details on the alternative options provided for each decision are given in the
279 following subsections. Note that some combinations of modeling options were clearly incompatible
280 with one another (options C1 and C2, for instance, cannot work with option D2). As a result, these
281 combinations were removed from the framework.

282 Another important feature of this modular framework is the systematic smoothing of all model
283 thresholds using infinitely differentiable approximants, as recommended by Kavetski and Kuczera
284 (2007) and Fenicia et al. (2011). The purpose here is twofold: first, to facilitate the calibration process
285 by removing any unnecessary (and potentially detrimental) discontinuities from the gradients of the
286 objective functions; and second, to provide a more realistic description of hydrological processes
287 across the catchment ([Moore and Clarke, 1981](#); Moore, 2007).

288
289

3.1.1. Snow accumulation and melt (decisions A and B)

290 Snow accumulation and melt components deal with the representation of snow processes at the
291 catchment scale. All modeling options rely on a single conceptual store to accumulate snow during the
292 winter months and release water during the melt season. Decision A refers to the partitioning of
293 precipitation into rain, snow or a mixture of rain and snow. Decision B refers to the representation of
294 snowmelt processes. Option A1 is the only hypothesis implemented to evaluate the relative abundance
295 of rain and snow. A logistic distribution is used in this option instead of usual temperature thresholds
296 to implicitly account for spatial variations in rain/snow partitioning over the catchment. In contrast,
297 three modeling options drawing upon the temperature-index approach (Hock, 2003) are available for
298 the evaluation of snowmelt rates (options B1a, B1b, B1c). Option B1a relies on a constant melt factor
299 while options B1b and B1c allow for temporal variability in the melt factor to reflect seasonal changes
300 in the energy available for melt. A recent example of option B1c can be found in [Clark et al.](#)
301 ([20092005](#)). Option B1b has been previously applied by [Schreider et al. \(1997\)](#) but at the grid cell
302 scale. Finally, it is worth noting that a smoothing kernel proposed by (Kavetski and Kuczera, 2007)
303 was introduced in the state equation of the snow reservoir to ignore residual snow remaining in the
304 reservoir outside the snowmelt season ([see Eq. \(1\)](#)).

305
306

3.1.2. Runoff generation (decision C)

307 Runoff generation components determine how much of a rainfall or snowmelt event is
308 available for runoff, lost through evapotranspiration or temporarily stored in soils and surface
309 depressions. Many models rely on a conceptual store to keep track of the catchment moisture status
310 and generate runoff as a function of both current and antecedent precipitation. Here, an assortment of
311 four commonly used methods is available. Option C1 is the only one in which no moisture accounting
312 store is required to estimate the contributing rainfall or snowmelt (see Fig. [32](#)). Actual
313 evapotranspiration then represents the only process involved in the production of runoff from
314 precipitation or snowmelt. The remaining options make use of moisture accounting stores and
315 distribution functions (see Table 1) to estimate the proportion of the basin generating runoff. An

316 important distinction is made between option C2, in which runoff generation occurs only during
317 rainfall or snowmelt events, and option C3, in which a leakage from the moisture accounting store
318 remains possible even after rainfall or snowmelt has ceased. Examples of these two moisture
319 accounting options can be found, respectively, in the HBV (e.g. Seibert and Vis, 2012) and PDM
320 (Moore, 2007) rainfall-runoff models. Alternative distribution functions are available in the literature,
321 for instance in the GR4j (Perrin *et al.*, 2003) and FLEX (Fenicia *et al.*, ~~2008b~~2006) models, but the
322 rationale behind their use remains the same. Actual evapotranspiration is computed from the estimated
323 PE using either a constant coefficient (option C1) or a function of the catchment moisture status
324 (options C2 and C3).

325 *3.1.3. Runoff transformation and routing (decisions D to F)*

326 Runoff transformation components account for all the retention and translation processes
327 occurring as water moves through the catchment. In practice, junction elements (decision D) and lag
328 functions (decision E) are typically combined with one or more conceptual stores (decision F) to
329 represent the effects of different flow pathways on the runoff process (both timing and volume).
330 Additional elements in the form of lag functions or conceptual stores can also be used to reflect water
331 routing in the channel network. However, in this study channel routing elements were considered
332 useless at a daily time step. All the modeling options available for decision F consist of two stores.
333 These can be arranged in parallel (options F1a and F1b), in series (options F2a and F2b), or in a
334 combination of both (options F3a and F3b). In each case, one of the stores has a nonlinear behavior
335 while the other reacts linearly. Two types of nonlinear response are provided: one that relies on
336 smoothed thresholds and different storage coefficients (options F1b, F2b and F3b), and the other that
337 relies on power laws (options F1a, F2a and F3a). Options F1a and F1b are based on the classical
338 parallel transfer function used in many conceptual models, such as the PDM (Moore, 2007) and
339 IHACRES (Jakeman *et al.*, 1993) models, where one store stands for a relatively quick catchment
340 response and the other for a slower response. The structure of options F3a and F3b is very close to the
341 response routine of the HBV model (e.g. Seibert and Vis, 2012). Note that some combinations of
342 modeling options were deemed unacceptable and thus not considered (e.g. D3–E1–F1a or D3–E1–
343 F1b).

344

345

346 3.2. Multi-objective optimization

347

348

3.2.1. Principle

349 In optimization problems with at least two conflicting objectives, a set of solutions rather than
350 a unique one exists because of the trade-offs between these objectives. A Pareto-optimal solution is
351 achieved when it cannot be improved upon without degrading at least one of its objective criteria. The
352 set of Pareto-optimal solutions for a given model is often called the “Pareto set” and the set of criteria
353 corresponding to this Pareto set is usually referred to as the “Pareto front”.

354

355

3.2.2. The NSGA-II algorithm

356 The Non-dominated Sorted Genetic Algorithm II (NSGA-II) (Deb, 2002) was selected to
357 calibrate the models implemented within the multiple-hypothesis framework. This algorithm has been
358 used successfully in a number of recent hydrological studies (see e.g. Khu and Madsen, 2005; Bekele
359 and Nicklow, 2007; De Vos and Rientjes, 2007; Fenicia *et al.*, 2008a; Shafii and De Smedt, 2009) and
360 has the advantage of not needing any additional parameter (other than those common to all genetic
361 algorithms, i.e. the initial population and the number of generations). Its most distinctive features are
362 the use of a binary tournament selection, a simulated binary crossover and a polynomial mutation
363 operator. For brevity’s sake, the detailed instructions of the algorithm and the conditions of its
364 application to rainfall-runoff modeling cannot be discussed further here. Instead, the reader is referred
365 to the aforementioned literature.

366

367 *3.2.3. Simulation periods and assessment criteria*

368 The simulation period was divided into a rather dry calibration period (1997–2011) and a
 369 relatively humid validation period (1982–1996). These two periods were chosen based on data
 370 availability to represent contrasted climate conditions: the two periods are separated by a shift in the
 371 IPO index, as explained in Sect [2.3.1.3.2.1](#).

372 Four criteria were chosen to evaluate the models built within the multiple-hypothesis
 373 framework. The first three of them are common to both calibration and validation periods while the
 374 fourth criterion differs between the two.

375 The first criterion (NSE) is related to the estimation of high flows and draws upon the Nash-
 376 Sutcliffe Efficiency metric:

$$\text{Crit1} = 1 - \text{NSE} = \sum_{d=1}^N (Q_{\text{obs}}^d - Q_{\text{sim}}^d)^2 / \sum_{d=1}^N (Q_{\text{obs}}^d - \overline{Q_{\text{obs}}})^2 \quad (2)$$

377 Where Q_{obs}^d and Q_{sim}^d are the observed and simulated discharges for day d , and N is the number of
 378 days with available observations.

379 The second criterion (NSE_{\log}) is related to the estimation of low flows and draws upon a modified, log
 380 version of the first criterion:

$$\text{Crit2} = 1 - \text{NSE}_{\log} = \sum_{d=1}^N (\log(Q_{\text{obs}}^d) - \log(Q_{\text{sim}}^d))^2 / \sum_{d=1}^N (\log(Q_{\text{obs}}^d) - \log(\overline{Q_{\text{obs}}}))^2 \quad (3)$$

381 The third criterion quantifies the mean annual volume error (VE_M) made in the estimation of the water
 382 balance of the catchment:

$$\text{Crit3} = \text{VE}_M = \sum_{y=1}^{N_{\text{years}}} (|V_{\text{obs}}^y - V_{\text{sim}}^y| / V_{\text{obs}}^y) / N_{\text{years}} \quad (4)$$

383 Where V_{obs}^y and V_{sim}^y are the observed and simulated volumes for year y , and N_{years} is the number of
 384 years of the simulation period.

385 The fourth criterion (Crit4) differs between the two simulation periods. In calibration, snow-covered
 386 areas (SCA) estimated from the MODIS data were used to evaluate the consistency of snow-
 387 accounting modeling options in terms of snow presence or absence ~~in the basin at the catchment scale~~.
 388 The objective was to quantify the error made in simulating the seasonal dynamics of snow
 389 accumulation, storage and melt processes. Following Parajka and Blöschl (2008), the snow error (SE)
 390 was defined as the total number of days when the snow-accounting store of options B1a, B1b and B1c
 391 disagreed with the MODIS data as to whether snow was present in the basin (Fig. [43](#)). The number of
 392 days with simulation errors is eventually divided by the total number of days with available MODIS
 393 data to express SE as a percentage.

394 In validation, a cumulated volume error was used to replace the snow error criterion that could not be
 395 computed due to a lack of remotely-sensed data over this period:

$$\text{Crit4} = \text{VE}_C = \left| \sum_{y=1}^{N_{\text{years}}} V_{\text{obs}}^y - \sum_{y=1}^{N_{\text{years}}} V_{\text{sim}}^y \right| / \sum_{y=1}^{N_{\text{years}}} V_{\text{obs}}^y \quad (5)$$

397

398 3.3. Model selection, model analysis and ensemble modeling

399 Finally, a total of 72 model structures were implemented and tested within the multi-objective and
 400 multiple-hypothesis frameworks. In addition to their names and for purposes of simplicity, these 72

401 model hypotheses are given a number from 1 to 72 corresponding to their order of appearance in the
402 simulation process (see e.g. Sect 4.1.).

403 Model hypotheses can be thought of as points x in the space of performance measures. One
404 possible way to locate these points in space is to consider that each coordinate $(x_i)_{i=1..4}$ of x is given
405 by the best performance obtained along the Pareto front of model x with respect to the i^{th} criterion
406 described in Sect 3.3.2. A clustering technique based on the fuzzy c-means algorithm (Bezdek et al.,
407 1983) and the initialization procedure developed by Chiu (1994) was chosen to explore this multi-
408 objective space and identify natural groupings among model hypotheses. To facilitate comparison
409 between calibration and validation, the clustering operations were repeated independently for each
410 period. The whole experiment, from model building to multi-objective optimization and cluster
411 identification, was repeated several times to ensure that the final composition of the clusters remains
412 the same.

413 Once the composition of each cluster was established, it was possible to identify a set of ‘best-
414 performing’ clusters for each simulation period, i.e. a set of clusters with the smallest Euclidian
415 distances to the origin of the objective space. The model structures of these ‘best-performing’ clusters
416 can be regarded as equally acceptable representations of the system. An important indicator of
417 structural uncertainty is the extent to which the simulation bounds derived from the Pareto sets of
418 these models reproduce the various features of the observed hydrograph. The overall uncertainty
419 envelope should be wide enough to include ~~most a large proportion~~ of the observed discharge but not
420 so wide that its representation of the various aspects of the hydrograph (rising limb, peak discharge,
421 falling limb, baseflow) becomes meaningless. ~~In general, one will seek to reduce as much as possible
422 the width of the envelope while maximizing the number of observations enclosed within the bounds.~~
423 In this study, priority was given to maintaining at its lowest value the number of outlying observations
424 before searching for the best combination of models which minimized the envelope area. This was
425 achieved iteratively through the following steps:

- 426
- 427 1. Start with an initial ensemble composed of the N_{max} models identified as members of the
428 best-performing clusters in both calibration and validation (i.e. models which fail the
429 validation test are ruled out).
- 430 2. From now on, consider only the calibration period.
431 Add up the N_{max} individual simulation envelopes that can be obtained from the Pareto sets of
432 the N_{max} models (hereafter referred to as the ‘Pareto-envelopes’).
- 433 3. Estimate the maximum number of observations enclosed within the resulting overall envelope,
434 $N_{obs}(N_{max})$, and calculate the area of this envelope, $Area(N_{max})$.
- 435 4. For $k = 1$ to N_{max}
 - 436 a. Identify the $\binom{N_{max}}{N_{max} - k}$ possible combinations of N_{max} models taken $N_{max} - k$ at a time.
 - 437 b. For each of these combinations
 - 438 - Add up the individual Pareto-envelopes of the $N_{max} - k$ models and calculate the
439 number of observations enclosed within the bounds of the resulting overall envelope,
440 $N_{obs}(N_{max} - k)$.
 - 441 - If $N_{obs}(N_{max} - k) = N_{obs}(N_{max})$
442 If $Area(N_{max} - k) < Area(N_{max} - k + 1)$
443 Accept the current combination.
 - 444 If $N_{obs}(N_{max} - k) < N_{obs}(N_{max})$
445 Reject the current combination.
 - 446 c. If all the possible combinations of $N_{max} - k$ models are rejected, break the loop. The final
447 ensemble of models to consider is the last accepted combination of $N_{max} - k + 1$ models.
 - 448

449 4. RESULTS

450

451 4.1. Model hypotheses evaluation

452

453 4.1.1. Cluster analysis

454 The 72 model hypotheses can be grouped into 5 clusters in calibration and 6 in validation. Table 3
455 displays the coordinates of the cluster centroids and gives, for each cluster, the number of points with
456 membership values above 50%. Figure 4-5 shows the projections of these clusters onto three possible
457 two-dimensional (2D) subspaces of the objective space (the three other subspaces being omitted for
458 brevity's sake). Each cluster is given a rank (from 1 to 5 or 6) reflecting its distance from the origin of
459 the coordinate system. As is evident from both Fig. 4-5 and Table 3, most of the best-performing
460 structures can be found in Cluster 1. This is particularly clear in the planes defined by the high-flow
461 (Crit1) and low-flow (Crit2) criteria (Figure 4-5), where all clusters tend to line up along a diagonal
462 axis (dashed line). In contrast, a small trade-off between Cluster 1 and Cluster 2 can be observed in
463 calibration in the plane defined by the high-flow (Crit1) and volume error (Crit3) criteria: models from
464 Cluster 2 (respectively Cluster 1) tend to perform slightly better than those from Cluster 1
465 (respectively Cluster 2) with respect to Crit3 (respectively Crit1). However, this trade-off disappears
466 in validation. Similar comments can be made about the other 2D subspaces (not shown here). In the
467 following analysis, Cluster 1 will be considered as the only best-performing cluster. This cluster
468 encompasses 24 members in calibration as against 15 in validation, indicating that several model
469 structures do not pass the validation test (namely models no. 30, 32, 49, 52, 53, 55, 66, 67, 69 and 72,
470 as shown in Table 4).

471 Several observations can be made regarding the composition of Cluster 1 in both simulation
472 periods. As can be seen from the values listed in Table 4, it is not possible to pick out a single,
473 unambiguous model hypothesis that would perform better than the others with respect to all criteria.
474 On the one hand, there appears to be several equally acceptable structures for each individual criterion.
475 Models no. 22 (A1-B1a-C3-D2-E1-F2b), 46 (A1-B1b-C3-D2-E1-F2b) and 54 (A1-B1c-C1-D3-
476 E2-F1b), for instance, yield very similar values of the high-flow criterion (Crit1), despite huge some
477 differences in their modeling options. This illustrates the equifinality of model structures in
478 reproducing one aspect of the system behavior. On the other hand, some structures seem more
479 appropriate to the simulation of high flows or snow dynamics while others appear to be better at
480 reproducing low flows or estimating the annual water balance of the catchment. This indicates trade-
481 offs between model structures in reproducing several aspects of the system behavior. It is however
482 possible to identify some recurring patterns among the modeling options present in (or absent from)
483 Cluster 1 in both periods. First, option B1c is the most represented snowmelt-accounting hypothesis,
484 despite an increase in the number of alternative options (B1a, B1b) in validation. More strikingly,
485 option C2 is totally absent from Cluster 1 in both periods. Single-flux combinations (C1-D1 and C3-
486 D2) and their splitting counterparts (C1-D3 and C3-D1) tend to be equally well-represented, thus
487 providing evidence of significant equifinality among these conceptual representations. Finally, runoff
488 transformation options based on a threshold-like behavior (F1b, F2b and F3b) account for 75% of
489 model hypotheses in calibration and over 90% in validation. In particular, option F3a turns out to be
490 completely absent from Cluster 1 in both periods while models based on option F2a (no. 49, 55, 67
491 and 69) fail the validation test. On the opposite, option F2b is particularly well-represented.

492

493

4.1.2. Pareto analysis

494 In general, valuable insight can be gained from the mapping of Pareto fronts in the space of
495 performance measures. While a full description of all the Pareto fronts obtained in calibration is not
496 possible here due to space limitations, two emblematic model hypotheses are used to illustrate this
497 point. Figure 5-6 shows the Pareto-optimal solutions of models no. 49 (A1-B1c-C1-D1-E1-F2a) and
498 50 (A1-B1c-C1-D1-E1-F2b) plotted in two dimensions for different combinations of two of the four
499 objective functions used in calibration. Note that these two models differ only in their runoff
500 transformation options (F2a vs. F2b) so that the comparison can be made in a controlled way. Trade-
501 offs between the high-flow (Crit1) and low-flow (Crit2) criteria are clearly more important with option
502 F2a (Fig. 5a6a) than with option F2b (Fig. 5b6b). This means that option F2a is less efficient in
503 reproducing simultaneously high and low flows and explains why this option disappears from Cluster
504 1 in validation. By contrast, the other pairs of criteria (Crit1-Crit3, Crit1-Crit4) displayed in Fig. 5-6
505 appear to be less useful in differentiating between the two models.

506 Further insight into the structural strengths and weaknesses of model hypotheses can be
507 obtained by determining how parameter values vary along the Pareto fronts of the models. A large

508 'Pareto range' in some parameters indicates structural deficiencies in the corresponding model
509 components (see e.g. Gupta *et al.*, 1998) or a lower sensitivity of model outputs to those parameters
510 (Engeland *et al.*, 2006). For purposes of clarity, Fig. 6-7 focuses on eight illustrative structures
511 identified as members of Cluster 1 in calibration. The models are paired in such a way that two models
512 of the same pair differ in only one modeling option. Thus, the effects of potential interactions between
513 model constituents are more likely to be detected. Parameter values are ~~normalised~~ normalized using
514 the lower and upper limits given in Table 2 so that all of them lie between 0 and 1. Different colors are
515 used to indicate the parameter sets associated with the smallest high-flow (in black), low-flow (in red),
516 volume (in blue) and snow (in green) errors. To what extent these colored solutions converge toward
517 the same parameter values or diverge from each other determines the level of parameter identifiability
518 of each model hypothesis. As regards snow-accounting options, a distinction can be made between
519 snow accumulation parameters (T_S and m_S), whose ranges of variation appear to be large in all cases,
520 and snowmelt parameters (T_M , f_M , r_1 , r_2 , f_1 , f_2), whose levels of identifiability depend on interactions
521 with the other model components. In Fig. 6a7a, the Pareto range of snowmelt parameters decreases in
522 width when moving from option B1a to B1b and using the combination of options C3–D2–E1. Yet
523 changing this combination into C3–D1–E2 has the opposite effect (Fig. 6b7b): parameter uncertainty
524 now decreases when moving from option B1b to B1a. As regards runoff transformation parameters (α ,
525 N_b , K_2 , K_3 , δ , S_C and K_4), the black and red solutions are closer to each other when options F2b (Fig.
526 6a7a, 6b7b and 6e7c) and F1b (Fig. 6d7d) are used. By contrast, options F2a (Fig. 6e7c) and F1a (Fig.
527 6d7d) require very different parameter sets to adequately simulate both low and high flows. Again,
528 this suggests that runoff transformation options based on a threshold-like behavior may be more
529 consistent with the observed data than those based on a power law relationship. It should be noted,
530 however, that relatively large Pareto ranges in some runoff transformation parameters (e.g. K_2 and K_3)
531 may still be required to obtain small volume and snow errors at the same time as high low-flow and
532 high-flow performances (e.g. models no. 44 and 54). Interestingly, the black, red and blue solutions of
533 models no. 49, 50, 53 and 54 also converge towards the same low values of parameter K_C
534 (evapotranspiration coefficient) independently of runoff transformation options.

535 Drawing any conclusion at this stage about the links between parameter identifiability and model
536 performance might be somewhat hazardous. Other examples (not shown here) show that a model
537 structure may have highly identifiable parameter values in calibration and yet not be suited to the
538 conditions prevailing in validation. Also, a reduction of parameter uncertainty as is the case with
539 options F2b and F1b often comes with a greater number of parameters.

540 Finally, a better understanding of the reasons why some models, or modeling options, work
541 better than others is provided by the simulation bounds (or Pareto-envelopes) derived from the Pareto
542 sets of these models. Figure 7-8 shows the Pareto-envelopes of the SWE internal state variable
543 obtained with three competing model hypotheses (no. 6, 30 and 54) differing only in their snowmelt-
544 accounting options (respectively B1a, B1b and B1c). Note that only the last two of these models (30,
545 54) belong to Cluster 1 in calibration (see Table 4). Simulated snow accumulation starts later than
546 expected with all modeling options (B1a, B1b and B1c). As will be further discussed in Sect 5.2., this
547 is likely to indicate systematic errors in the input precipitation and/or MODIS-based SCA data. On the
548 whole, the envelope widths suggest a reduction in the uncertainty associated with the prediction of
549 snow seasonal dynamics when moving from option B1a to option B1c. This is consistent with the
550 mean annual snow errors reported in Table 4, which are significantly lower with option B1c
551 independently of the other model options. It must be acknowledged, however, that even this option
552 (B1c) fails to capture the seasonal dynamics of snow accumulation and melt during several years of
553 the period. The release of water from the snow-accounting store of model no. 54 continues well after
554 the end of the observed snowmelt season in 2008, 2009, 2010 and 2011. On the contrary, the
555 simulated snowmelt season tends to end sooner than expected with model no. 30 in 2003, 2004, 2005
556 and 2006. In that case, options B1b and B1c appear to be somewhat complementary.

557 ~~4.2. Comparison with the physical features of the catchment~~

558 ~~4.2.1. Snow accumulation and melt~~

559

The relatively large Pareto bounds obtained for parameters T_s and m_s with nearly all model hypotheses indicate that mixed conditions of rain and snow are likely to occur across a large range of temperatures. This may be due to the lumped representation of the snow accumulation process and the necessity to implicitly account for spatial variations in rain/snow partitioning across the catchment. Likewise, the relatively high values of parameter K_c (> 0.2) obtained with the green solutions (smallest snow errors) of models no. 50, 53 and 54 (Fig. 6) might indicate a need to compensate for the absence of sublimation scheme in the available snow modeling options. The sine function used in option B1c appears to be better suited to the estimation of the melt factor than the other options tested in this study (B1a, B1b). The degree-day method implemented in option B1a has a physical basis (Ohmura, 2011). Yet some components in the energy balance of snow covered areas cannot be fully captured by temperature alone nor easily reduced to a simple formula (Hoek, 2003). In semi-arid central Andes (29–30°S), small zenith angles and a thin, dry and cloud-free atmosphere during most of the year make incoming shortwave radiation the most important source of seasonal variations in the energy available for melt (see e.g. Aberman *et al.*, 2013). As a result, the seasonal timing of snowmelt is expected to show greater year-to-year stability, which may explain the relative success of option B1c when compared to option B1b.

4.2.2. Runoff generation

The absence of option C2 in Cluster 1 in both simulation periods suggests that moisture accounting components may not be essential to the conceptual modeling of this semi-arid Andean catchment. Most of the land cover is, indeed, dominated by barren to sparsely vegetated exposed rocks, boulders and rubble with poor soil development outside the valleys. This setting may also explain the relatively low values of parameter K_c obtained with the black, red and blue solutions shown in Fig. 6.

4.2.3. Runoff transformation and routing

The high representation of options F2a and F2b in Cluster 1 suggests that the catchment actually behaves as a serial system and may reveal a better correspondence with its overall physical structure. The overall organization of fluxes in the catchment, from high elevations toward the valleys and then northward to the outlet, can be conceptualized as a series of two hydraulically connected reservoirs: one standing for the mountain blocks (upstream reservoir) and the other for the alluvial valleys (downstream reservoir). Of course, this interpretation needs to be qualified, since other runoff transformation options (F1a, F1b and F3b) have proved to yield equally acceptable simulations despite significant differences in their model structures.

4.3.4.2. Representation of structural uncertainties

This Section deals with the identification and use of an ensemble of equally acceptable model structures to quantify and represent the uncertainty arising from the system non-identifiability. Figure 97 shows the overall uncertainty envelope obtained with the 8 model structures whose combination minimizes the envelope area in calibration while holding constant the number of outlying observations (see Sect 3.3.). Over 82% of discharge observations are captured by the envelope in both simulation periods. Interestingly, this number exceeds the best N_{par} value obtained in calibration with the individual Pareto-envelopes (see Table 4), which shows how necessary it is to consider an ensemble of model structures. In validation, however, a better combination could be identified since several models of Cluster 1 display significantly higher N_{par} values (Table 4). On the whole, the comparison of the observed hydrograph with the simulation bounds of the envelope shows a good match of rising limbs and peak discharges in both simulation periods, but a less accurate fit of falling limbs during at least one major (in 1987–88) and two minor (in 2005–06 and 2007–08) events. The slower recession of the observed hydrograph might indicate a delayed contribution of one or more catchment compartments that cannot be described by any of the modeling options available in the multiple-hypothesis framework.

609
610
611
612

5. DISCUSSION & CONCLUSION

613 This study aimed at reducing structural uncertainty in the modeling of a semi-arid Andean catchment
614 where lumped conceptual models remain largely under-used. To overcome the current lack of
615 information on model adequacy in this catchment, a modular modeling framework (MMF) relying on
616 six model-building decisions was developed to generate 72 competing model structures. Four
617 assessment criteria were then chosen to calibrate and evaluate these models over a 30-year period
618 using the concept of Pareto-optimality. This strategy was designed to characterize both the parameter
619 uncertainty arising from each model's structural deficiencies (i.e. model inadequacy) and the
620 ambiguity associated with the choice of model components (i.e. model non-uniqueness). Finally, a
621 clustering approach was taken to identify natural groupings in the multi-objective space. Overall, the
622 greatest source of uncertainty was found in the connection between runoff generation and runoff
623 transformation components (decisions D and E). However, the results also showed a significant drop
624 in the number of plausible representations of the system. After validation, 14 model structures among
625 the 24 identified in calibration as the best-performing ones were finally considered as equally
626 acceptable.

627 Interestingly, both rejected and accepted hypotheses appeared closely related to particular types of
628 snowmelt-accounting (decision B), runoff generation (decision C) and runoff transformation (decision
629 D) modeling options, suggesting possible links to some physical features of the catchment. For
630 instance, the frequent occurrence of option C1 and the absence of option C2 among the set of best-
631 performing structures indicate that moisture-accounting components may not be essential to the
632 conceptual modeling of this catchment. Most of the land cover is, indeed, dominated by barren to
633 sparsely vegetated exposed rocks, boulders and rubble with poor soil development outside the valleys.
634 This setting may also explain the relatively low values of parameter K_C obtained with the black, red
635 and blue solutions shown in Fig. 6. Likewise, the frequency of options F2a and F2b in the best-
636 performing cluster suggests that the catchment actually behaves as a 'serial' system. The overall
637 organization of fluxes in the catchment, from high elevations toward the valleys and then northward to
638 the outlet, can be conceptualized as a series of two hydraulically connected reservoirs: one standing
639 for the granitic mountain blocks (upstream reservoir) and the other for the alluvial valleys
640 (downstream reservoir). Similar results were also obtained for smaller catchments in Luxembourg
641 characterized by relatively impervious bedrocks and lateral water flows (Fenicia et al., 2014). The
642 results also provided some evidence of a strong threshold behavior at the catchment scale (options
643 F1b, F2b and F3b) compared to the smoother power laws of options F1a, F2a and F3a. However,
644 further research would be needed to track the origin of this behavior, which might be related at some
645 point to connectivity levels in the fractured and till-mantled areas of the mountain blocks. As regards
646 snowmelt, the frequent occurrence of option B1c in the best-performing cluster in calibration may
647 indicate a need to account for processes which the degree-day method implemented in option B1a does
648 not fully capture. In semi-arid central Andes (29–30°S), small zenith angles and a thin, dry and cloud-
649 free atmosphere during most of the year make incoming shortwave radiation the most important
650 source of seasonal variations in the energy available for melt (e.g. Pellicciotti et al., 2008; Abermann
651 et al., 2013). While this dominant source of energy cannot be accounted for by temperature alone, the
652 seasonal timing of snowmelt is also expected to show a greater year-to-year stability, which may
653 explain the relative success of option B1c when compared to option B1b. Of course, these
654 hypothesized relationships between some physical characteristics of the catchment and specific
655 modeling options need to be further qualified. Differentiating between physically adequate and purely
656 numerical solutions will always seem somewhat hazardous in the case of lumped conceptual models.
657 For instance, a small number of models among those identified as the best-performing ones also rely
658 on parallel (F1a, F1b) and intermediate (F3b) runoff transformation options. Also, the relative
659 proportions of snowmelt-accounting options B1a, B1b and B1c, appears much more balanced in
660 validation, where no snow error criterion could be applied, than in calibration. Although this was not
661 our objective in this paper, comparative studies including several similar or contrasted catchments
662 would be required to better understand how different model structures relate to different physical

663 settings. Such understanding is of primary importance to the choice of conceptual models in climate
664 change impact studies.

665 Another important issue related to model identification is the extent to which the 'principle of
666 parsimony' can be applied to differentiate between a large number of model hypotheses. Many authors
667 rightly consider that a maximum of 5 to 6 parameters should be accepted in calibration when using a
668 single objective function. Efstratiadis and Koutsoyiannis (2010) extended this empirical rule to the
669 case of multi-objective schemes by allowing « a ratio of about 1:5 to 1:6 between the number of
670 criteria and the number of parameters to optimize ». For a multi-objective scheme based on four
671 criteria (as in the present study), this leads to consider 20 to 24-parameter models as still being
672 parsimonious. This will certainly seem unreasonable to many modelers because, as Efstratiadis and
673 Koutsoyiannis (2010) also pointed out, the various criteria used are generally not independent of each
674 other. In our case, for instance, the information added by the low-flow criterion may not be so
675 different from that already introduced by the high-flow criterion. By contrast, the snow criterion tends
676 to add new information on the snow-related parameters. From this perspective, it is noteworthy that
677 most rejected hypotheses among the 24 identified in calibration as members of Cluster 1 had more
678 than 11 free parameters, with only one having 9 parameters. The principle of parsimony, however,
679 cannot be used to further discriminate between the remaining 14 best-performing hypotheses. For
680 instance, model no. 54 (12 parameters) performs better than model no. 2 (9 parameters) with respect to
681 the high-flow criterion.

682 Eventually, the number of models used to represent structural uncertainty was reduced by
683 searching for which minimal set of models maximized the number of observations covered by the
684 ensemble of Pareto-envelopes. It is important to make clear that model inadequacy and non-
685 uniqueness were evaluated here in non-probabilistic terms. In particular, the Pareto-envelopes derived
686 for each model structure quantify only the uncertainty arising from the trade-offs between competing
687 criteria and do not have a predefined statistical meaning (Engeland et al., 2006). Consequently, the
688 overall simulation bounds shown in Figure 8 cannot be easily interpreted as 'confidence bands'.
689 Although discussing the adequacy of non-probabilistic approaches to structural uncertainty was far
690 beyond the scope of this study, it is interesting to analyze the reasons why between 15% and 20% of
691 the observations remained outside the overall simulated envelope in both calibration and validation.
692 To a large extent, this lack of performance can be attributed either to an insufficient coverage of the
693 hypothesis and objective spaces or to uncertainties in the precipitation and streamflow data that were
694 overlooked in this study.

695 First, the choice of Pareto-optimality to characterize structural uncertainty can be criticized for
696 leading to the rejection of many behavioral parameter sets (i.e. being close to, but not part of, the
697 Pareto front) that might have been Pareto-optimal with different performance measures, calibration
698 data or input errors (e.g. Freer et al., 2003; Beven, 2006). Also, this concept should not be confused
699 with that of equifinality. Both notions agree that it is not possible to identify a single, best solution to
700 the calibration problem and that multiple parameters sets should be retained to give a proper account
701 of model uncertainty. However, the Pareto set of solutions represents the minimum parameter
702 uncertainty that can be achieved when several criteria are considered simultaneously with no *a priori*
703 preference for one over the others (Gupta et al., 2003). By contrast, two parameter sets are said to be
704 equifinal (in a statistical sense) if they can be regarded as equally acceptable with respect to a given
705 model outcome. For a proper assessment of parameter equifinality, more probabilistic approaches
706 should be taken (Madsen, 2000; Huisman *et al.*, 2010). In the context of multiple-hypothesis testing, a
707 meticulous selection of the assessment criteria is also critical to avoid rejecting some modeling options
708 for the wrong reasons. For instance, the snow error criterion was shown to have a great influence on
709 the identification of snow-accounting components, as much more ambiguity between the various
710 available options was observed during the validation period when this criterion could not be used.
711 Also, like any other multiple-hypothesis framework, the MMF developed in this study suffers from an
712 insufficient coverage of the hypothesis space (Gupta et al., 2012). The parameterization of
713 evapotranspiration, for example, was not considered as an independent model-building decision. Only
714 one formula was applied to calculate potential evapotranspiration and the possibility to retrieve actual
715 evapotranspiration from downstream water stores was not provided. Likewise, the runoff
716 transformation process was described using only two water stores, of which only one was assumed to
717 have a nonlinear behavior. Future work to improve the conceptual modeling of the Claro River

718 catchment should include the testing of new or refined hypotheses to allow for the use of additional
719 auxiliary data (e.g. observed snow heights, irrigation water-use).

720 More fundamentally, our ability to discriminate among the competing model hypotheses was
721 constrained by inevitable errors in the input and output data sets. In particular, the comparison of
722 simulated SWE levels and MODIS-based SCA estimates revealed some uncertainty in the estimation
723 of precipitation inputs and confirmed previous results obtained by Favier et al. (2009). Some
724 precipitation events occurring in the early winter may not be captured by the gauging network (< 3200
725 m a.s.l.) used for the interpolation of precipitation across the catchment. These errors may add to
726 systematic volume errors caused by wind, wetting and evaporation losses at the gauge level, leading to
727 an overall underestimation of precipitation, as indicated by the rough estimate of the catchment-scale
728 water balance given in Sect 2. It was also possible to highlight some errors in the streamflow data. The
729 observed streamflow was ‘naturalized’ by simply adding back the estimated historical water
730 abstractions (Sect. 2.2). When applied on a daily basis, this process inevitably adds some uncertainty
731 to streamflow values because a significant part of surface-water abstractions actually return to the river
732 system within a few days due to conveyance and field losses. In general, ignoring these return flows
733 would lead to overestimating daily natural flows. In this paper, however, the actual water withdrawals
734 were not known with precision but only as percentages of the nominal water rights – these percentages
735 being fixed on a monthly basis by the authorities to account for variations in water availability. The
736 combined impact of streamflow and precipitation errors on the assessment of structural uncertainty
737 thus remained unknown. Further research is currently underway to integrate the effects of water
738 abstractions and crop water-use in the hydrological modeling process (Hublart et al., 2015; see also
739 Kiptala et al., 2014 for another approach). From a multiple-hypothesis perspective, the modeling of
740 irrigation water-use should be regarded as a testable model component in its own right.

741 ~~This study provided an opportunity to combine a modular modeling approach with a multi-criteria~~
742 ~~evaluation scheme to reduce structural uncertainty in the conceptual modeling of a large Andean~~
743 ~~catchment over a 30-year period. In particular, it demonstrated the benefits of using the concept of~~
744 ~~Pareto efficiency to discriminate among several competing model structures. Among the 72~~
745 ~~hypotheses tested, the results showed that 58 model hypotheses can be rejected as inappropriate.~~
746 ~~However, 14 other hypotheses were shown to yield equally acceptable representations of the~~
747 ~~catchment hydrological functioning in both calibration and validation. Further, the simulation~~
748 ~~envelopes derived from the Pareto sets of 8 model structures among the 14 best performing ones were~~
749 ~~used to represent the *minimum* structural uncertainty that could be obtained with this modeling~~
750 ~~framework. The rejection of some hypotheses was closely related to particular types of model~~
751 ~~components or modeling options. For instance, option C2, in which runoff generation requires the~~
752 ~~filling a moisture accounting store, can be ruled out from the set of plausible runoff generation~~
753 ~~representations. It is noteworthy that most rejected hypotheses among the 24 identified in calibration~~
754 ~~as the best performing ones have more than 11 free parameters, with only one rejected hypothesis~~
755 ~~having 9 parameters. Thus, more parsimonious models seem to better withstand changes in the climate~~
756 ~~conditions. The principle of parsimony, however, cannot be used to further discriminate between the~~
757 ~~remaining best performing hypotheses. For instance, model no. 54 (12 parameters) performs better~~
758 ~~than model no. 2 (9 parameters) with respect to the high-flow criterion.~~

759 ~~There remains several ways to improve this assessment of structural uncertainty and model~~
760 ~~suitability. In particular, the concept of Pareto optimality should not be confused with that of~~
761 ~~equifinality. Of course, both notions agree that it is not possible to identify a single, best solution to~~
762 ~~the calibration problem and that multiple parameter sets should be retained to give a proper account of~~
763 ~~model uncertainty. However, the Pareto set of solutions represents the minimum parameter uncertainty~~
764 ~~that can be achieved when several criteria are considered simultaneously with no *a priori* preference~~
765 ~~for one over the others (Gupta *et al.*, 2003). By contrast, two parameter sets are said to be equifinal if~~
766 ~~they can be regarded as equally acceptable in a statistical sense with respect to one particular criterion~~
767 ~~(for more details on these differences, see Engeland *et al.*, 2006). From this perspective, the choice of~~
768 ~~Pareto optimality to characterize model uncertainty can be criticized for leading to the rejection of~~
769 ~~many behavioural parameter sets (i.e. being close to, but not part of, the Pareto front) that might have~~
770 ~~been Pareto optimal with different performance measures, calibration data or errors in the input data~~
771 ~~(e.g. Freer *et al.*, 2003; Beven, 2006). One possible way to address this limitation and improve model~~
772 ~~transposability in time has been suggested by Gharari *et al.* (2013). The idea is to divide the~~

773 calibration period into k sub-periods and identify parameter sets (in the whole parameter space) which
774 minimize the distance to the k Pareto fronts of these sub-periods. For a proper assessment of parameter
775 equifinality, however, Bayesian frameworks should be considered (Madsen, 2000; Huisman *et al.*,
776 2010).

777 The use of Pareto envelopes to quantify structural uncertainty is also questionable in that it fails to
778 account for all discharge observations, as shown in Table 4. While this failure can be partly remedied
779 within a multiple hypothesis framework (MHF), Fig. 8 shows that the overall uncertainty envelope
780 obtained by merging the Pareto envelopes of 8 competing model hypotheses still leaves out a
781 significant part of the observations. Indeed, like any other modular framework, the MHF developed in
782 this study suffers from an insufficient coverage of the hypothesis space (Gupta *et al.*, 2012). The
783 parameterization of evapotranspiration, for example, was not considered as an independent model-
784 building decision. Only one formula was applied to calculate potential evapotranspiration and the
785 possibility to retrieve actual evapotranspiration from downstream water stores was not provided.
786 Likewise, the runoff transformation process was described using only two water stores, of which only
787 one was assumed to have a nonlinear behavior. Future work to improve the conceptual modeling of
788 the Claro River Catchment should include the testing of new or refined hypotheses to allow for the use
789 of additional auxiliary data (e.g. groundwater levels). Competing alternatives to the lumped mode used
790 in this study should also be included within the MHF. For example, semi-lumped approaches in which
791 snow accumulation and melt components are applied at the grid-cell level provide an interesting way
792 to improve the use of snow cover data without increasing too much computational requirements. In
793 this way, catchment wide snow covered areas (SCA) can be simulated and directly compared to
794 MODIS based data. Daily rainfall and snowmelt amounts are then integrated over all grid cells to be
795 used as catchment averaged inputs in the subsequent spatially-lumped model components (see e.g.
796 Schreider *et al.*, 1997). This improved MHF should then be applied to other mesoscale catchments to
797 better understand how the specific features of each catchment relate to specific model requirements.
798 Such understanding is of primary importance for the use of conceptual models in climate change
799 impact studies.

800 Finally, our ability to discriminate among the competing model hypotheses was constrained by
801 inevitable errors in the input and output data sets. In particular, the comparison of simulated SWE
802 levels and MODIS-based SCA estimates revealed considerable uncertainty in the estimation of
803 precipitation inputs. Some precipitation events occurring in the early winter are not captured by the
804 gauging network (< 3000 m a.s.l.) used for the interpolation of precipitation across the catchment.
805 These errors add to the systematic volume errors caused by wind, wetting and evaporation losses at the
806 gauge level, leading to an overall underestimation of precipitation, as indicated by the rough
807 estimation of catchment scale water balance given in Sect 2. It was also possible to highlight some
808 errors in the streamflow data. Part of these errors might be associated with uncertainties in the
809 estimation of natural streamflow. Further research is therefore required to better integrate the effect of
810 water abstractions in the hydrological modeling process. From a multiple hypothesis perspective, the
811 modeling of irrigation water withdrawals should be regarded as a testable model component in its own
812 right.

814 **Acknowledgements** The authors are very grateful to the Centro de Estudios Avanzados en Zonas
815 Áridas (CEAZA) for its essential logistic support during the field missions and to Gustavo Freixas
816 from the *Dirección General de Agua* (Chile) for providing the necessary streamflow data. The authors
817 also thank S. Lhermitte, D. López and S. MacDonell for providing the MODIS data used in this study
818 and S. Gascoin for informal advice and much useful discussion. Moreover, the authors thank the two
819 anonymous reviewers for their interest to this work and for their useful comments that helped to
820 improve the article.

822 REFERENCES

823 Abermann, J., Kinnard, C., and MacDonell, S.: Albedo variations and the impact of clouds on glaciers in the
824 Chilean semi-arid Andes, *J. Glaciol.*, 60, 183–191, 2013.

- 825 | ~~Anderson, P. W.: More is different, *Science*, 177, 393–396, 1972.~~
- 826 | Bekele, E. G. and Nicklow, J. W.: Multi-objective automatic calibration of SWAT using NSGA-II, *J. Hydrol.*,
827 341, 165–176, 2007.
- 828 | Beven, K.: Prophecy, reality and uncertainty in distributed hydrological modelling, *Adv. Water Resour.*, 16, 41–
829 51, 1993.
- 830 | Beven, K.: A Manifesto for the Equifinality Thesis, *J. Hydrol.*, 320, 18–36, 2006.
- 831 | ~~Bezdek~~Bezdek, J. C., Ehrlich, R., and Full, W.: FCM: The fuzzy c-means clustering algorithm, *Comput. Geosci.*,
832 10, 191–203, 1983.
- 833 | Birkel, C., Tetzlaff, D., Dunn, S. M., and Soulsby, C.: Towards a simple dynamic process conceptualization in
834 rainfall–runoff models using multi-criteria calibration and tracers in temperate, upland catchments. *Hydrol.*
835 *Process.*, 24, 260–275, doi: 10.1002/hyp.7478, 2010.
- 836 | Blöschl, G. and A. Montanari: Climate change impacts–throwing the dice?, *Hydrol. Process.*, 24, 374–381, 2010.
- 837 | ~~Blöschl, G. and M. Sivapalan: Scale issues in hydrological modelling: a review, *Hydrol. Processes*, 9, 251–290,
838 1995.~~
- 839 | Boyle, D. P., Gupta, H. V., and Sorooshian, S.: Toward improved calibration of hydrologic models: Combining
840 the strengths of manual and automatic methods, *Water Resour. Res.*, 36, 3663–3674,
841 doi:10.1029/2000WR900207, 2000.
- 842 | Buytaert, W. and K. Beven: Models as multiple working hypotheses: hydrological simulation of tropical alpine
843 wetlands, *Hydrol. Process.*, 25, 1784–1799, 2011.
- 844 | Capell, R., Tetzlaff, D., and Soulsby, C.: Can time domain and source area tracers reduce uncertainty in
845 rainfall-runoff models in larger heterogeneous catchments?, *Water Resour. Res.*, 48, W09544,
846 doi:10.1029/2011WR011543, 2012.
- 847 | Caviedes, C. N. and Paskoff, R.: Quaternary glaciations in the Andes of north-central Chile, *J. Glaciol.*, 14, 155–
848 169, 1975.
- 849 | Centro del Agua para Zonas Áridas y semiáridas de América Latina y el Caribe (CAZALAC): Aplicación de
850 metodologías para determinar la eficiencia de uso del agua – Estudio de caso en la Región de Coquimbo.
851 Informe Técnico, Gobierno Regional, Santiago (Chile), 2006.
- 852 | ~~Centro de Información de Recursos Naturales (CIREN): Descripción de suelos, materiales y símbolos, IV
853 Región. Publ. CIREN N°129, Estudio Agrológico, Chile, 2005 (in Spanish).~~
- 854 | ~~Centro de Información de Recursos Naturales (CIREN), Catastro Frutícola, IV Región de Coquimbo (Chile),
855 Publ. CIREN, 2011.~~
- 856 | Chiu, S.: Fuzzy model identification based on cluster estimation, *J. Intell. Fuzzy Syst.*, 2, 267–278, 1994.
- 857 | Clark, M. P., Slater, A. G., Barrett, A. P., Hay, L. E., McCabe, G. J., Rajagopalan, B., and Leavesley, G. H.:
858 Assimilation of snow covered area information into hydrologic and landsurface models, *Adv. Water*
859 *Resour.*, 29, 1209–1221, 2006.
- 860 | Clark, M. P., Slater, A. G., Rupp, D. E., Woods, R. A., Vrugt, J. A., Gupta, H. V., Wagener, T., and Hay, L. E.:
861 Framework for Understanding Structural Errors (FUSE): A modular framework to diagnose differences
862 between hydrological models, *Water Resour. Res.*, 44, W00B02, doi:10.1029/2007WR006735, 2008.
- 863 | Clark, M., Hreinsson, E. O., Martinez, G., Tait, A., Slater, A., Hendrikx, J., Owens, I., Gupta, H., Schmidt, J., and
864 Woods, R.: Simulations of seasonal snow for the South Island, New Zealand, *J. Hydrol.*, 48, 41–58, 2009.
- 865 | Clark, M. P., Kavetski, D., and Fenicia, F.: Pursuing the method of multiple working hypotheses for hydrological
866 modeling, *Water Resour. Res.*, 47, W09301, doi:10.1029/2010WR009827, 2011.
- 867 | Collet, L., Ruelland, D., Borrell-Estupina, V., Dezetter, A., and Servat, E.: Integrated modelling to assess long-
868 term water supply capacity of a meso-scale Mediterranean catchment, *Sci. Total Environ.*, 461–462, 528–
869 540, 2013.
- 870 | Coxon, G., Freer, J., Wagener, T., Odoni, N. A., and Clark, M. P.: Diagnostic evaluation of multiple hypotheses
871 of hydrological behaviour in a limits-of-acceptability framework for 24 UK catchments, *Hydrol. Process.*,
872 doi:10.1002/hyp.10096, online first, 2013.

- 873 De Vos, N. J., and Rientjes, T. H. M.: Multi-objective performance comparison of an artificial neural network
874 and a conceptual rainfall-runoff model, *Hydrolog. Sci. J.*, 52, 397–413, doi: 10.1623/hysj.52.3.397, 2007.
- 875 Deb, K., Pratap, A., Agarwal, S., and Meyarivan, T.: A fast and elitist multiobjective genetic algorithm: NSGA-II,
876 *IEEE T. Evolut. Comput.*, 6, 181–197, 2002.
- 877 Dooge, J.: Looking for hydrologic laws, *Water Resour. Res.*, 22, 46S–58S, doi:10.1029/WR022i09Sp0046S,
878 1986.
- 879 Dooge, J.: Searching for Simplicity in Hydrology, *Surv. Geophys.*, 18, 511–534, 1997.
- 880 Efstratiadis, A., and Koutsoyiannis, D.: One decade of multi-objective calibration approaches in hydrological
881 modelling: a review, *Hydrolog. Sci. J.*, 55, 58–78, 2010.
- 882 Ehret, U., Gupta, H. V., Sivapalan, M., Weijjs, S. V., Schymanski, S. J., Blöschl, G., Gelfan, A. N., Harman, C.,
883 Kleidon, A., Bogaard, T. A., Wang, D., Wagener, T., Scherer, U., Zehe, E., Bierkens, M. F. P., Di
884 Baldassarre, G., Parajka, J., van Beek, L. P. H., van Griensven, A., Westhoff, M. C., and Winsemius, H. C.:
885 Advancing catchment hydrology to deal with predictions under change, *Hydrol. Earth Syst. Sci.*, 18, 649–
886 671, doi:10.5194/hess-18-649-2014, 2014.
- 887 Engeland, K., Braud, I., Gottschalk, L., and Leblois, E.: Multi-objective regional modelling, *J. Hydrol.*, 327,
888 339–351, 2006.
- 889 Favier, V., Falvey, M., Rabatel, A., Praderio, E., and López, D.: Interpreting discrepancies between discharge and
890 precipitation in high-altitude area of Chile's Norte Chico region (26–32°S), *Water Resour. Res.*, 45,
891 W02424, doi:10.1029/2008WR006802, 2009.
- 892 ~~Fenicia, F., Solomatine, D. P., Savenije, H. H. G., and Matgen, P.: Soft combination of local models in a multi-
893 objective framework, *Hydrol. Earth Syst. Sci.*, 11, 1797–1809, 2007.~~
- 894 Fenicia, F., McDonnell, J. J., and Savenije, H. H. G.: Learning from model improvement: On the contribution of
895 complementary data to process understanding, *Water Resour. Res.*, 44, W06419,
896 doi:10.1029/2007WR006386, ~~2008~~2008a.
- 897 ~~Fenicia, F., Savenije, H. H. G., Matgen, P. and Pfister, L.: Understanding catchment behavior through stepwise
898 model concept improvement, *Water Resour. Res.*, 44, W01402, doi:10.1029/2006WR005563, 2008b.~~
- 899 Fenicia, F., Kavetski, D., and Savenije, H. H. G.: Elements of a flexible approach for conceptual hydrological
900 modeling: 1. Motivation and theoretical development, *Water Resour. Res.*, 47, W11510,
901 doi:10.1029/2010WR010174, 2011.
- 902 Fenicia, F., Kavetski, D., Savenije, H. H. G., Clark, M. P., Schoups, G., Pfister, L., and Freer, J.: Catchment
903 properties, function, and conceptual model representation: is there a correspondence?, *Hydrol. Process.*,
904 28, 2451–2467, doi: 10.1002/hyp.9726, 2014.
- 905 ~~Fowler, H. J., and Kilsby, C. G.: Using regional climate model data to simulate historical and future river flows
906 in northwest England, *Climatic Change*, 80, 337–367, 2007.~~
- 907 Freer, J., Beven, K., and Peters, N.: Multivariate Seasonal Period Model Rejection Within the Generalised
908 Likelihood Uncertainty Estimation Procedure, in *Calibration of Watershed Models* (eds Q. Duan, H. V.
909 Gupta, S. Sorooshian, A. N. Rousseau and R. Turcotte), American Geophysical Union, Washington, D. C.,
910 69–87doi: 10.1029/WS006p0069, ~~2003~~2013.
- 911 ~~Gharari, S., Hrachowitz, M., Fenicia, F., and Savenije, H. H. G.: An approach to identify time-consistent model
912 parameters: sub-period calibration, *Hydrol. Earth Syst. Sci.*, 17, 149–161, doi:10.5194/hess-17-149-2013,
913 2013.~~
- 914 Gupta, H. V., Sorooshian, S., and Yapo, P. O.: Toward improved calibration of hydrologic models: Multiple and
915 noncommensurable measures of information, *Water Resour. Res.*, 34, 751–763, 1998.
- 916 Gupta, H. V., Bastidas, L. A., Vrugt, J. A., and Sorooshian, S.: Multiple criteria global optimization for watershed
917 model calibration, *Water Sci. Appl.*, 6, 125–132, 2003.
- 918 Gupta, H. V., Clark, M. P., Vrugt, J. A., Abramowitz, G., and Ye, M.: Towards a comprehensive assessment of
919 model structural adequacy, *Water Resour. Res.*, 48, W08301, doi:10.1029/2011WR011044, 2012.
- 920 Hock, R.: Temperature index melt modelling in mountain areas, *J. Hydrol.*, 282, 104–115, 2003.
- 921 ~~Horton, P., Schaefli, B., Mezghani, A., Hingray, B., and Musy, A.: Assessment of climate change impacts on
922 alpine discharge regimes with climate model uncertainty, *Hydrol. Processes*, 20, 2091–2109, 2006.~~

- 923 [Hrachowitz, M., Savenije, H. H. G., Blöschl, G., McDonnell, J. J., Sivapalan, M., Pomeroy, J. W., Arheimer, B.,](#)
 924 [Blume, T., Clark, M. P., Ehret, U., Fenicia, F., Freer, J. E., Gelfan, A., Gupta, H. V., Hughes, D. A., Hut, R.](#)
 925 [W., Montanari, A., Pande, S., Tetzlaff, D., Troch, P. A., Uhlenbrook, S., Wagener, T., Winsemius, H. C.,](#)
 926 [Woods, R. A., Zehe, E., and Cudennee, C.: A decade of Predictions in Ungauged Basins \(PUB\) — a review,](#)
 927 [Hydrol. Sci. J., 58, 1198–1255, DOI:10.1080/02626667.2013.803183, 2013.](#)
- 928 Hublart, P., Ruelland, D., Dezetter, A., and Jourde, H.: Modeling current and future trends in water availability
 929 for agriculture on a semi-arid and mountainous Chilean catchment, in: Cold and Mountain Region
 930 Hydrological Systems Under Climate Change: Towards Improved Projections, IAHS-AISH P., 360, 26–32,
 931 2013.
- 932 Hublart, P., Ruelland, D., Dezetter, A., and Jourde, H.: Assessing the capacity to meet irrigation water needs for
 933 viticulture under climate variability in the Chilean Andes, in: Hydrology in a Changing World:
 934 Environmental and Human Dimensions, Proc. 7th FRIEND Int. Conf., Montpellier, France, 24–28 February
 935 2014, IAHS-AISH P., 363, 209–214, 2014.
- 936 [Hublart, P., Ruelland, D., García de Cortázar Atauri, I., and Ibacache, A.: Assessing the reliability of conceptual](#)
 937 [hydrological modeling in a cultivated, drought-prone catchment of the Chilean Andes, in: Hydrologic Non-](#)
 938 [Stationarity and Extrapolating Models to Predict the Future, IAHS-AISH P. \(in press\), 2015.](#)
- 939 Huisman, J. A., Rings, J., Vrugt, J. A., Sorg, J., Vereecken, H.: Hydraulic properties of a model dike from
 940 coupled Bayesian and multi-criteria hydrogeophysical inversion, *J. Hydrol.*, 380, 62–73, 2010.
- 941 ~~[Instituto Nacional de Estadísticas \(INE\): Catastro Vitícola Nacional 2007–2008, Publ. Anual, 2009.](#)~~
- 942 IPCC: Full Report: the Physical Science Basis, in: Contribution of Working Group I to the Fifth Assessment
 943 Report of the Intergovernmental Panel on Climate Change, Climate Change 2013, edited by: Stocker, T. F.,
 944 Qin, D., Plattner, G.-K., Tignor, M., Allen, S. K., Boschung, J., Nauels, A., Xia, Y., Bex, V., and Midgley, P.
 945 M., Cambridge University Press, Cambridge, UK and New York, NY, USA, 1261–1264, 2013.
- 946 Jakeman, A. J. and Hornberger, G. M.: How much complexity is warranted in a rainfall-runoff model?, *Water*
 947 *Resour. Res.*, 29, 2637–2649, 1993.
- 948 Jakeman, A. J., and Letcher, R. A.: Integrated assessment and modelling: features, principles and examples for
 949 catchment management, *Environ. Modell. Softw.*, 18, 491–501, 2003.
- 950 Jothityangkoon, C., Sivapalan, M., and Farmer, D. L.: Process controls of water balance variability in a large
 951 semi-arid catchment: downward approach to hydrological model development, *J. Hydrol.*, 254, 174–198,
 952 2001.
- 953 Jourde, H., Rochette, R., Blanc, M., Brisset, N., Ruelland, D., Freixas, G., and Oyarzun, R.: Relative
 954 contribution of groundwater and surface water fluxes in response to climate variability of a mountainous
 955 catchment in the Chilean Andes, in: Cold Regions Hydrology in a Changing Climate, IAHS-AISH P., 346,
 956 180–188, 2011.
- 957 Kalthoff, N., Fiebig-Wittmaack, M., Meißner, C., Kohler, M., Uriarte, M., Bischoff-Gauß, I., and Gonzales, E.:
 958 The energy balance, evapo-transpiration and nocturnal dew deposition of an arid valley in the Andes, *J.*
 959 *Arid Environ.*, 65, 420–443, 2006.
- 960 Kavetski, D., and Kuczera, G.: Model smoothing strategies to remove microscale discontinuities and spurious
 961 secondary optima in objective functions in hydrological calibration, *Water Resour. Res.*, 43, W03411, ,
 962 doi:10.1029/2006WR005195, 2007.
- 963 Kavetski, D., and Fenicia, F.: Elements of a flexible approach for conceptual hydrological modeling: 2.
 964 Application and experimental insights, *Water Resour. Res.*, 47, W11511, doi:10.1029/2011WR010748,
 965 2011.
- 966 Khu, S. T., and Madsen, H.: Multiobjective calibration with Pareto preference ordering: An application to
 967 rainfall-runoff model calibration, *Water Resour. Res.*, 41, W03004, 10.1029/2004WR003041, 2005.
- 968 [Kiptala, J. K., Mul, M. L., Mohamed, Y. A. and van der Zaag, P.: Modelling stream flow and quantifying blue](#)
 969 [water using a modified STREAM model for a heterogeneous, highly utilized and data-scarce river basin in](#)
 970 [Africa, Hydrol. Earth Syst. Sci., 18, 2287–2303, 2014.](#)
- 971 Kirchner, J. W.: Getting the right answers for the right reasons: Linking measurements, analyses, and models to
 972 advance the science of hydrology, *Water Resour. Res.*, 42, WR004362, doi:10.1029/2005WR004362, 2006.
- 973 Kokkonen, T. S., and Jakeman, A. J.: A comparison of metric and conceptual approaches in rainfall-runoff
 974 modeling and its implications, *Water Resour. Res.*, 37, 2345–2352, 2001.

- 975 Krueger, T., Freer, J., Quinton, J. N., Macleod, C. J. A., Bilotta, G. S., Brazier, R. E., Butler, P., and Haygarth, P.
976 M.: Ensemble evaluation of hydrological model hypotheses, *Water Resour. Res.*, 46, W07516.
977 doi:10.1029/2009WR007845, 2010.
- ~~978 Lang, H., and Braun, L.: On the information content of air temperature in the context of snow melt estimation.
979 In: Molnar, L., (Ed.), Hydrology of Mountainous Areas, Proceedings of the Strbske Pleso Symposium
980 1990: IAHS Publ. no. 190, 347–354, 1990.~~
- ~~981 Leavesley, G. H., Markstrom, S. L., Restrepo, P. J., and Viger, R. J.: A modular approach to addressing model
982 design, scale, and parameter estimation issues in distributed hydrological modelling, *Hydrol. Process.*, 16,
983 173–187, 2002.~~
- 984 Lee, G., Tachikawa, Y., and Takara, K.: Comparison of model structural uncertainty using a multi-objective
985 optimization method, *Hydrol. Process.*, 25, 2642–2653, 2011.
- ~~986 Loukas, A., Vasiliades, L., and Dalezios, N. R.: Climatic impacts on the runoff generation processes in British
987 Columbia, Canada, *Hydrol. Earth Syst. Sci.*, 6, 211–227, 2002.~~
- 988 Madsen, H.: Automatic calibration of a conceptual rainfall–runoff model using multiple objectives, *J. Hydrol.*,
989 235, 276–288, 2000.
- 990 MacDonell, S., Kinnard, C., Mölg, T., Nicholson, L., and Abermann, J.: Meteorological drivers of ablation
991 processes on a cold glacier in the semiarid Andes of Chile, *The Cryosphere*, 7, 1833–1870, doi:10.5194/tc-
992 7-1513-2013, 2013.
- 993 McDonnell, J. J., Sivapalan, M., Vaché, K., Dunn, S., Grant, G., Haggerty, R., Hinz, C., Hooper, R., Kirchner, J.,
994 Roderick, M. L., Selker, J., and Weiler, M.: Moving beyond heterogeneity and process complexity: A new
995 vision for watershed hydrology, *Water Resour. Res.*, 43, W07301, doi:10.1029/2006WR005467, 2007.
- 996 McMillan, H.: Effect of spatial variability and seasonality in soil moisture on drainage thresholds and fluxes in a
997 conceptual hydrological model, *Hydrol. Process.*, 26, 2838–2844, doi: 10.1002/hyp.9396, 2012a.
- 998 McMillan, H., Tetzlaff, D., Clark, M., and Soulsby, C.: Do time-variable tracers aid the evaluation of
999 hydrological model structure? A multimodel approach, *Water Resour. Res.*, 48, W05501,
1000 doi:10.1029/2011WR011688, 2012b.
- 1001 Michaud, J., and Sorooshian, S.: Comparison of simple versus complex distributed runoff models on a semi-arid
1002 watershed, *Water Resour. Res.*, 30, 593–605, 1994.
- 1003 Milano, M., Ruelland, D., Dezetter, A., Fabre, J., Ardoin-Bardin, S., and Servat, E.: Modeling the current and
1004 future capacity of water resources to meet water demands in the Ebro basin, *J. Hydrol.*, 500, 114–126,
1005 2013.
- 1006 Minville, M., Brissette, F., and Leconte, R.: Uncertainty of the impact of climate change on the hydrology of a
1007 nordic watershed, *J. Hydrol.*, 358, 70–83, 2008.
- 1008 Montecinos, A. and ~~Accituno, P.~~ Patrieio, A.: Seasonality of the ENSO-Related Rainfall Variability in Central
1009 Chile and Associated Circulation Anomalies, *J. Climate*, 16, 281–296, 2003.
- 1010 Moore, R. J.: The PDM rainfall-runoff model, *Hydrol. Earth Syst. Sci.*, 11, 483–499, 2007.
- ~~1011 Moore, R. J., and Clarke, R. T.: A distribution function approach to rainfall runoff modeling, *Water Resour. Res.*,
1012 17(5), 1367–1382, 1981.~~
- ~~1013 Olsson, J. A., and Andersson, L.: Possibilities and problems with the use of models as a communication tool in
1014 water ressource management, *Water Resour. Manage.*, 21, 97–110, 2007.~~
- ~~1015 Ohmura, A.: Physical basis for the temperature-based melt index method, *J. Appl. Meteorol.*, 40, 753–761, 2001.~~
- 1016 Oudin, L., Hervieu, F., Michel, C., Perrin, C., Andréassian, V., Anctil, F., and Loumagne, C.: Which potential
1017 evapotranspiration input for a lumped rainfall–runoff model?: Part 2—Towards a simple and efficient
1018 potential evapotranspiration model for rainfall–runoff modelling, *J. Hydrol.*, 303, 290–306, 2005.
- 1019 Parajka, J., and Blöschl, G.: The value of MODIS snow cover data in validating and calibrating conceptual
1020 hydrologic models, *J. Hydrol.*, 358, 240–258, 2008.
- ~~1021 Pellicciotti, F., Helbing, J., Rivera, A., Favier, V., Corripio, J., Araos, J., Sicart, J.-E. and Carenzo, M.: A study of
1022 the energy balance and melt regime on Juncal Norte Glacier, semi-arid Andes of central Chile, using melt
1023 models of different complexity, *Hydrol. Process.*, 22, 3980–3997. doi: 10.1002/hyp.7085, 2008.~~

- 1024 Perrin, C., Michel, C., and Andréassian, V.: Improvement of a parsimonious model for streamflow simulation, *J.*
1025 *Hydrol.*, 279, 275–289, 2003.
- 1026 Pourrier, J., Jourde, H., Kinnard, C., Gascoïn, S., and Monnier, S.: Glacier meltwater flow paths and storage in a
1027 geomorphologically complex glacial foreland: The case of the Tapado glacier, dry Andes of Chile (30° S), *J.*
1028 *Hydrol.*, 519, 1068–1083, doi:10.1016/j.jhydrol.2014.08.023, 2014.
- 1029 Quintana, J. M. and Aceituno, P.: Changes in the rainfall regime along the extratropical west coast of South
1030 America (Chile): 30–43°S, *Atmósfera*, 25, 1–22, 2012.
- 1031 Refsgaard, J. C., and Knudsen, J.: Operational validation and intercomparison of different types of hydrological
1032 models, *Water Resour. Res.*, 32, 2189–2202, 1996.
- 1033 Ruelland, D., Brisset, N., Jourde, H., and Oyarzun, R.: Modelling the impact of climatic variability on the
1034 groundwater and surface flows from a mountainous catchment in the Chilean Andes, in: *Cold Regions*
1035 *Hydrology in a Changing Climate*, IAHS-AISH P., 346, 171–179, 2011.
- 1036 Ruelland, D., Ardoin-Bardin, S., Collet, L., and Roucou, P.: Simulating future trends in hydrological regime of a
1037 large Sudano-Sahelian catchment under climate change, *J. Hydrol.*, 424–425, 207–216, 2012.
- 1038 Ruelland, D., Dezetter, A., and Hublart, P.: Sensitivity analysis of hydrological modelling to climate forcing in a
1039 semi-arid mountainous catchment, in: *Hydrology in a Changing World: Environmental and Human*
1040 *Dimensions*, Proc. 7th FRIEND Int. Conf., Montpellier, France, 24–28 February 2014, IAHS-AISH P., 363,
1041 145–150, 2014.
- 1042 ~~Rutllant, J. and Fuenzalida, H.: Synoptic aspects of the central Chile rainfall variability associated with the~~
1043 ~~Southern Oscillation, *Int. J. Climatol.*, 11, 63–76, 1991.~~
- 1044 Savenije, H. H. G.: HESS Opinions "The art of hydrology", *Hydrol. Earth Syst. Sci.*, 13, 157–161,
1045 doi:10.5194/hess-13-157-2009, 2009.
- 1046 Schaeffli, B., Harman, C. J., Sivapalan, M., and Schymanski, S. J.: HESS Opinions: Hydrologic predictions in a
1047 changing environment: behavioral modeling, *Hydrol. Earth Syst. Sci.*, 15, 635–646, doi:10.5194/hess-15-
1048 635-2011, 2011.
- 1049 Schreider, S., Whetton, P. H., Jakeman, A. J., and Pittock, A. B.: Runoff modelling for snow-affected catchments
1050 in the Australian alpine region, eastern Victoria, *J. Hydrol.*, 200, 1–23, 1997.
- 1051 ~~Schulz, N., Boisier, J. P., and Aceituno, P.: Climate change along the arid coast of northern Chile, *Int. J.*
1052 ~~*Climatol.*, 32, 1803–1814, 2011.~~~~
- 1053 Seibert, J.: Multi-criteria calibration of a conceptual runoff model using a genetic algorithm, *Hydrol. Earth Syst.*
1054 *Sci.*, 4, 215–224, doi:10.5194/hess-4-215-2000, 2000.
- 1055 Seibert, J. and McDonnell, J. J.: On the dialog between experimentalist and modeler in catchment hydrology:
1056 Use of soft data for multicriteria model calibration, *Water Resour. Res.*, 38, W01241, doi:
1057 10.1029/2001WR000978, 2002.
- 1058 Seibert, J. and Vis, M. J. P.: Teaching hydrological modeling with a user-friendly catchment-runoff-model
1059 software package, *Hydrol. Earth Syst. Sci.*, 16, 3315–3325, 2012.
- 1060 Shafii, M. and De Smedt, F.: Multi-objective calibration of a distributed hydrological model (WetSpa) using a
1061 genetic algorithm, *Hydrol. Earth Syst. Sci.*, 13, 2137–2149, 2009.
- 1062 ~~Sivapalan, M., Blöschl, G., Zhang, L., and Vertessy, R.: Downward approach to hydrological prediction, *Hydrol.*
1063 ~~*Process.*, 17, 2101–2111, 2003.~~~~
- 1064 Sivapalan, M.: Pattern, process and function: elements of a unified theory of hydrology at the catchment scale,
1065 *Encyclopedia of Hydrological Sciences*, doi:10.1002/0470848944.hsa012, online first, 2006.
- 1066 ~~Sivapalan, M., Blöschl, G., Zhang, L., and Vertessy, R.: Downward approach to hydrological prediction, *Hydrol.*
1067 ~~*Process.*, 17, 2101–2111, 2003.~~~~
- 1068 Smith, T. J., and Marshall, L. A.: Exploring uncertainty and model predictive performance concepts via a
1069 modular snowmelt-runoff modeling framework, *Environ. Modell. Softw.*, 25, 691–701, 2010.
- 1070 Son, K., and Sivapalan, M.: Improving model structure and reducing parameter uncertainty in conceptual water
1071 balance models through the use of auxiliary data, *Water Resour. Res.*, 43, W01415,
1072 doi:10.1029/2006WR005032, 2007.

- 1073 Souvignet, M: Climate Change Impacts on Water Availability in the Semiarid Elqui Valley, Chile, Ph.D. thesis,
1074 Cologne University of Applied Sciences, Institute for Technology in the Tropics, 110 pp., 2007.
- 1075 Souvignet, M., Hartmut, G., Lars, R., Kretschmer, N., and Oyarzún, R.: Statistical downscaling of precipitation
1076 and temperature in north-central Chile: an assessment of possible climate change impacts in an arid Andean
1077 watershed, *Hydrol. Sci. J.*, 55, 41–57, 2010.
- 1078 Squeo, F. A., Veit, H., Arancio, G., Gutiérrez, J. R., Arroyo, M. T. K., and Olivares, N.: Spatial heterogeneity of
1079 high mountain vegetation in the Andean desert zone of Chile (30°S), *Mt. Res. Dev.*, 13, 203–209, 1993.
- 1080 Staudinger, M., Stahl, K., Seibert, J., Clark, M. P., and Tallaksen, L. M.: Comparison of hydrological model
1081 structures based on recession and low flow simulations, *Hydrol. Earth Syst. Sci.*, 15, 3447–3459,
1082 doi:10.5194/hess-15-3447-2011, 2011.
- 1083 Strauch, G., Oyarzun, J., Fiebig-Wittmaack, M., González, E., and Weise, S. M.: Contributions of the different
1084 water sources to the Elqui river runoff (northern Chile) evaluated by H/O isotopes, *Isot. Environ. Health S.*,
1085 42, 303–322, 2006.
- 1086 Verbist, K., Robertson, A. W., Cornelis, W. M., and Gabriels, D.: Seasonal predictability of daily rainfall
1087 characteristics in central northern Chile for dry-land management, *J. Appl. Meteorol. Clim.*, 49, 1938–1955,
1088 2010.
- 1089 Vicuña, S., Garreaud, R., and McPhee, J.: Climate change impacts on the hydrology of a snowmelt driven basin
1090 in semiarid Chile, *Climatic Change*, 105, 469–488, [20102011](#).
- 1091 Wagener, T., Lees, M. J., and Wheeler, H. S.: A toolkit for the development and applications of parsimonious
1092 hydrological models, in: *Mathematical Models of Large Watershed Hydrology*, vol. 1, edited by: Singh, V.
1093 P. and Frevert, D., Water Resources Publishers, Highland Ranch, CO, 87–136, 2002.
- 1094 ~~Wagener, T., Wheeler, H. S., and Gupta, H. V.: *Rainfall runoff modelling in gauged and ungauged catchments*.
1095 London: Imperial College Press, 332 p, 2004.~~
- 1096 ~~Wainwright, J., and Mulligan, M. (Eds): *Environmental modelling – Finding simplicity in complexity*.
1097 Chichester, John Wiley & Sons, Ltd., 2004.~~
- 1098 Xu, C.-Y., and Singh, V. P.: Review on regional water resources assessment models under stationary and
1099 changing climate, *Water Resour. Manage.*, 18, 591–612, 2004.
- 1100 Young, G., Zavala, H., Wandel, J., Smit, B., Salas, S., Jimenez, E., Fiebig, M., Espinoza, R., Diaz, H., and
1101 Cepeda, J.: Vulnerability and adaptation in a dryland community of the Elqui Valley, Chile, *Climatic
1102 Change*, 98, 245–276, 2010.
- 1103 ~~Zhang, X., Srinivasan, R., and Van Liew, M.: *On the use of multi algorithm, genetically adaptive multi objective
1104 method for multi site calibration of the SWAT model*, *Hydrol. Process.*, 24, 955–969, 2010.~~

1105 **TABLES & CAPTIONS**

1106

1107 **Table 1.** Constitutive equations of fluxes between the various components of the modeling options described in
 1108 Fig. 2. Parameter (in italic) significations and units are detailed in Table 2. P: catchment-averaged daily
 1109 precipitation; Rain: rain fraction of precipitation P; Snow: snow fraction of precipitation P; T: catchment-
 1110 averaged daily temperature; PE: catchment-averaged daily potential evapotranspiration; AE: catchment-averaged
 1111 daily actual evapotranspiration; $S_j, j \in [1,5]$: state variables of the conceptual stores; $Q_j, j \in [1,5]$: water fluxes
 1112 between the model components).

Options	Constitutive equations	Options	Constitutive equations
A1	$\text{Snow} = P / (1 + \exp[(T - T_s) / m_s])$ $= P(1 + \exp[(T - T_s) / m_s])$ $\text{Rain} = P - \text{Snow}$	C3	$Q_1 = (\text{Melt} + \text{Rain})[1 - (1 - S_1/S_m)^b]$ $Q_2 = K_1 S_1$
B1a, B1b, B1c	$\text{Melt} = MF(\bar{T} - \log[1 + \exp(-\bar{T})])$ with $\bar{T} = (T - T_M) / m_M$ and $m_M = 0.1^\circ\text{C}$	D1	$Q_3 = Q_2 \text{ and } Q_4 = Q_1$ or $Q_3 = Q_1$
B1a	$MF = f_M m_M$	D2	$Q_3 = Q_1 + Q_2$
B1b	$MF = r_1 + r_2 T_{30}$ with T_{30} the mean temperature of the last 30 days	D3	$Q_3 = (1 - \alpha)Q_1$ $Q_4 = \alpha Q_1$
B1c	$MF = f_1 + f_2 \sin(0.551\pi + 2\pi d/366)$	E1	$Q_{j,\text{lag}} = Q_2$ with $j \in \{3,4\}$
C1	$AE = \min(\text{Melt} + \text{Rain}, K_C PE)$	E2	$Q_{j,\text{lag}}(t) = \sum_{i=1}^{N_b} \omega(i) Q_j(t - i + 1)$ with $\omega(i) = \int_{i-1}^i 2udu / N_b^2$
C2, C3	$AE = PE \min(1, S_1/S_m)$	F1a, F2a, F3a	$Q_5 = K_2 S_2^{1+\delta}$ $Q_6 = K_3 S_3$
C1	$Q_1 = \text{Melt} + \text{Rain}$	F1b, F2b, F3b	$Q_5 = K_4 S_2 + K_2 (\bar{S}_2 - \log[1 + \exp(-\bar{S}_2)])$ $Q_6 = K_3 S_3$ with $\bar{S}_2 = (S_2 - S_C) / m_C$ and $m_C = 0.1 \text{ mm}^{-1}$
C2	$Q_1 = (\text{Melt} + \text{Rain})(S_1/S_m)^\beta$	F3a, F3b	$Q_6 = DS_2$

1113

1114 | **Table 2.** Parameters used in the various modeling options with their signification and initial sampling. (*) The
 1115 | possible values for K_C were limited to a maximum of 0.5 to reflect the extreme aridity of the catchment.
 1116 |

Parameter	Options	Signification	Units	Initial range
T_S	A1	Rain / snow partitioning temperature threshold	°C	-10 – 10
m_S	A1	Rain / snow partitioning smoothing parameter	–	0.01 – 3
T_M	B1a, B1b, B1c	Snowmelt temperature threshold	°C	-10 – 10
f_M	B1a	Constant melt factor	°C.mm ⁻¹	0 – 10
r_1	B1b	Coefficient for computation of the variable melt factor	°C.mm ⁻¹	1 – 5
r_2	B1b	Coefficient for computation of the variable melt factor	°C.mm ⁻¹	1 – 5
f_1	B1c	Coefficient for computation of the variable melt factor	°C.mm ⁻¹	1 – 5
f_2	B1c	Coefficient for computation of the variable melt factor	°C.mm ⁻¹	1 – 5
K_C	C1	Evapotranspiration coefficient	–	0.05 – 0.5 (*)
S_m	C2, C3	Maximum storage capacity of the moisture-accounting store	mm	10 – 100
β	C2	Shape parameter	–	0.1 – 3
b	C3	Shape parameter of Pareto distribution	–	0.1 – 3
K_1	C3	Infiltration coefficient	d ⁻¹	0.001 – 0.7
α	D3	Splitting parameter	–	0.1 – 0.9
N_b	E2	Number of time steps in the lag routine	–	1 – 6
K_2	F1a to F3b	Storage coefficient	d ⁻¹	0.01 – 0.99
K_3	F1a to F3b	Storage coefficient	d ⁻¹	0.001 – 0.01 (F1a, F1b, F3a, F3b) 0.001 – 0.10-01 (F2a, F2b)
δ	F1a, F2a, F3a	Power law parameter of the non-linear store in the runoff transformation module	–	0 – 1
S_c	F1b, F2b, F3b	Threshold parameter of the non-linear store in the runoff transformation module	mm	10 – 300
D	F3a, F3b	Recharge coefficient	d ⁻¹	0.001 – 0.5
K_4	F1b, F2b, F3b	Storage coefficient	d ⁻¹	0.001 – 0.01

1117

1118 **Table 3.** Coordinates of the cluster centroids in the four-dimensional (4D) space of performance measures. The
 1119 number of models with membership values > 50% ($N_{50\%}$) is given for each cluster.
 1120

Calibration period (1997–2011)					
Cluster no.	Crit1 (1-NSE)	Crit2 (1-NSE _{log})	Crit3 (VE _M) (%)	Crit4 (SE) (%)	N _{50%}
1	0.15	0.25	10	9	24
2	0.23	0.30	10	10	24
3	0.49	0.58	23	11	10
4	0.60	0.62	25	16	13
5	0.92	0.97	33	20	1

Validation period (1982–1996)					
Cluster no.	Crit1 (1-NSE)	Crit2 (1-NSE _{log})	Crit3 (VE _M) (%)	Crit4 (VE _C) (%)	N _{50%}
1	0.24	0.21	14	3	15
2	0.32	0.29	15	4	25
3	0.38	0.31	15	5	8
4	0.51	0.42	25	23	8
5	0.61	0.44	27	27	11
6	0.61	0.51	30	33	5

1121

1122 **Table 4.** Detailed composition of Clusters 1 in calibration and validation. The tables indicate the numbers and
 1123 the names of the models as well as their number of parameters NP. For each criterion only the best performance
 1124 value obtained along the Pareto front is given. N_{par} (%) represents the proportion of observations enclosed within
 1125 the simulation bounds of each Pareto set of solutions. Asterisks are used to indicate the models which are not in
 1126 the best-performing group (Cluster 1) either in calibration or in validation.
 1127

Calibration period (1997–2011)							
Model no.	Model name (options)	NP	NSE	NSE _{log}	VE _M (%)	SE (%)	N_{par} (%)
2	A1–B1a–C1–D1–E1–F2b	9	0.87	0.76	10.6	11.2	76.0
4	A1–B1a–C1–D1–E1–F3b	10	0.84	0.77	10.4	11.2	53.2
8	A1–B1a–C1–D3–E2–F2b	11	0.83	0.75	11.7	11.1	76.5
20	A1–B1a–C3–D1–E2–F2b	12	0.83	0.76	10.0	11.4	60.0
22	A1–B1a–C3–D2–E1–F2b	11	0.90	0.77	10.4	11.2	64.1
26	A1–B1b–C1–D1–E1–F2b	10	0.87	0.77	10.1	11.5	58.4
30 (*)	A1–B1b–C1–D3–E2–F1b	12	0.84	0.70	9.8	11.4	69.6
32 (*)	A1–B1b–C1–D3–E2–F2b	12	0.83	0.71	11.1	11.4	68.4
44	A1–B1b–C3–D1–E2–F2b	13	0.89	0.77	10.6	11.4	63.4
46	A1–B1b–C3–D2–E1–F2b	12	0.90	0.76	10.7	11.4	45.4
49 (*)	A1–B1c–C1–D1–E1–F2a	9	0.82	0.73	10.9	7.0	67.0
50	A1–B1c–C1–D1–E1–F2b	10	0.86	0.77	10.4	7.0	67.4
52 (*)	A1–B1c–C1–D1–E1–F3b	11	0.85	0.72	8.8	8.1	65.7
53 (*)	A1–B1c–C1–D3–E2–F1a	11	0.79	0.76	10.8	7.0	63.8
54	A1–B1c–C1–D3–E2–F1b	12	0.90	0.78	11.5	7.5	55.7
55 (*)	A1–B1c–C1–D3–E2–F2a	11	0.80	0.73	10.7	7.0	54.5
56	A1–B1c–C1–D3–E2–F2b	12	0.85	0.75	10.8	7.6	76.3
65	A1–B1c–C3–D1–E2–F1a	12	0.83	0.78	8.0	7.7	65.0
66 (*)	A1–B1c–C3–D1–E2–F1b	13	0.81	0.77	9.6	6.8	63.5
67 (*)	A1–B1c–C3–D1–E2–F2a	12	0.81	0.75	10.7	7.0	73.7
68	A1–B1c–C3–D1–E2–F2b	13	0.85	0.74	10.6	6.8	74.5
69 (*)	A1–B1c–C3–D2–E1–F2a	11	0.82	0.73	10.6	7.0	51.8
70	A1–B1c–C3–D2–E1–F2b	12	0.87	0.76	10.7	7.5	76.4
72 (*)	A1–B1c–C3–D2–E1–F3b	13	0.81	0.71	9.8	7.1	69.0

1128
1129

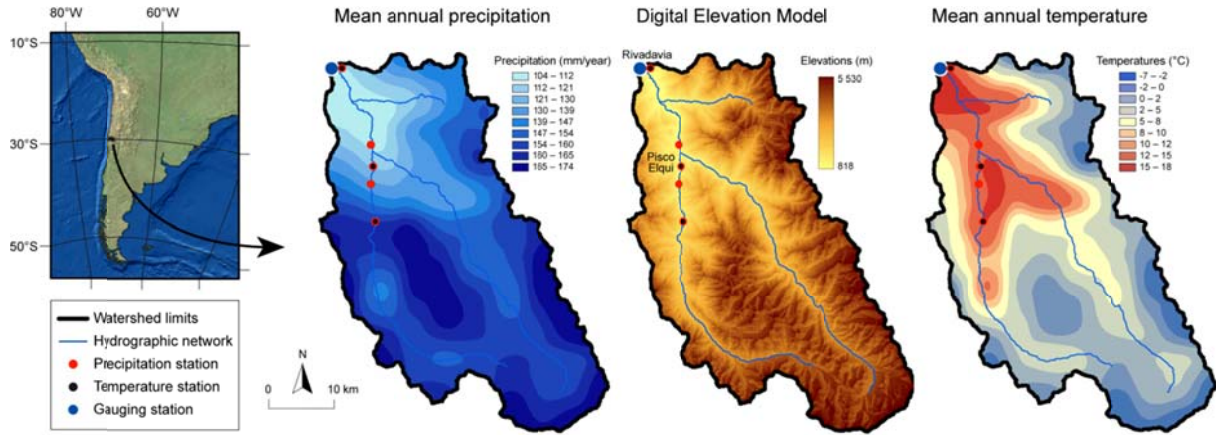
Validation period (1982–1996)							
Model no.	Model name	NP	NSE	NSE _{log}	VE _M (%)	VE _C (%)	N_{par} (%)
2	A1–B1a–C1–D1–E1–F2b	9	0.75	0.78	13.3	2.7	87.1
4	A1–B1a–C1–D1–E1–F3b	10	0.73	0.80	14.1	3.8	50.0
8	A1–B1a–C1–D3–E2–F2b	11	0.75	0.76	14.5	5.8	84.8
20	A1–B1a–C3–D1–E2–F2b	12	0.72	0.77	13.7	3.7	58.4
22	A1–B1a–C3–D2–E1–F2b	11	0.76	0.78	12.3	3.3	75.3
26	A1–B1b–C1–D1–E1–F2b	10	0.74	0.78	12.9	3.5	70.2
42 (*)	A1–B1b–C3–D1–E2–F1b	13	0.73	0.75	15.6	3.3	62.7
44	A1–B1b–C3–D1–E2–F2b	13	0.74	0.79	13.0	4.1	69.3
46	A1–B1b–C3–D2–E1–F2b	12	0.76	0.77	15.2	3.4	48.4
50	A1–B1c–C1–D1–E1–F2b	10	0.78	0.81	13.9	2.5	73.1
54	A1–B1c–C1–D3–E2–F1b	12	0.77	0.78	15.3	3.5	60.8
56	A1–B1c–C1–D3–E2–F2b	12	0.75	0.77	13.2	4.5	81.3
65	A1–B1c–C3–D1–E2–F1a	12	0.74	0.80	13.8	3.6	73.0
68	A1–B1c–C3–D1–E2–F2b	13	0.77	0.74	13.5	3.7	78.7
70	A1–B1c–C3–D2–E1–F2b	12	0.73	0.78	14.2	3.4	79.4

1130

1131 FIGURES & CAPTIONS

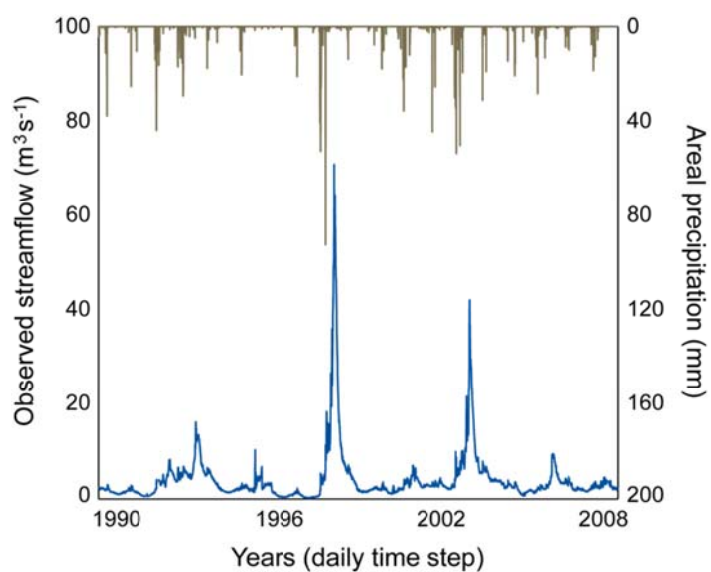
1132

1133 **Figure 1.** The Claro River Basin at Rivadavia (1515 km²) in Chile: topography and mean annual precipitation
1134 and temperature over 1982–2011 (based on Ruelland *et al.*, 2014). Several of the stations used in this study were
1135 located outside the catchment and therefore not displayed on the following maps.



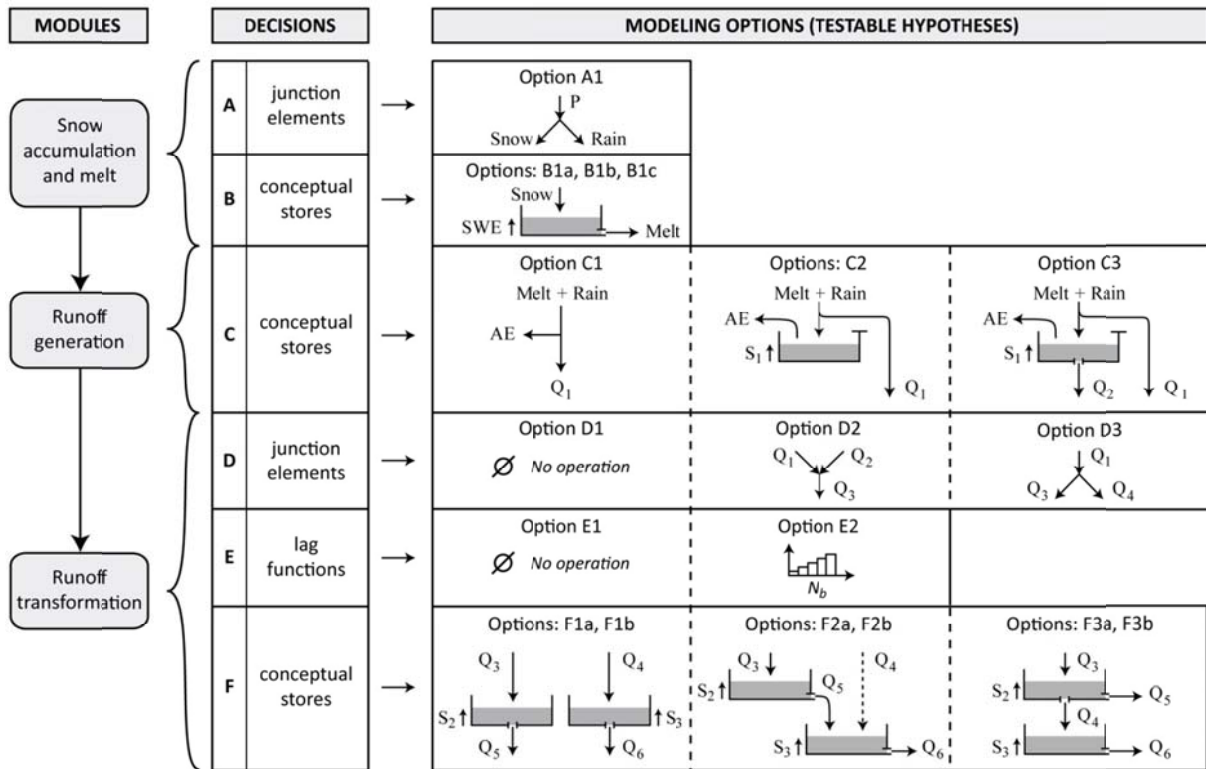
1136

1137 **Figure 2.** Interannual variability in precipitation and observed streamflow from 1989 to 2008. The hydrological
1138 year was defined from May to April so as to capture the snowmelt and peak flow seasons at mid-year.
1139 Streamflow values are those measured at the catchment outlet before accounting for water abstractions.
1140 Precipitation values are those obtained after interpolation.



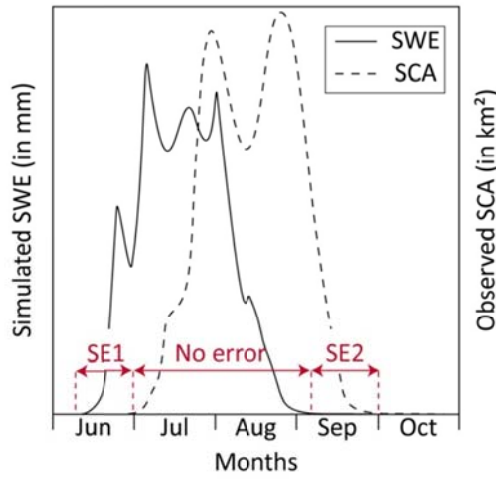
1141

1142 **Figure 3.** Overall architecture (modules), decision tree and available modeling options of the modular multiple-
 1143 hypothesis framework (P: catchment-averaged daily precipitation; SWE: snow water equivalent; AE: catchment-
 1144 averaged daily actual evapotranspiration; $S_j, j \in [1,5]$: state variables of the conceptual stores; $Q_j, j \in [1,5]$: water
 1145 fluxes between the model components).



1146

1147 **Figure 4.** Description of the snow error criterion. The overall snow error (SE) can be described as a sum of two
 1148 terms, SE1 and SE2, whose values are given by a confusion matrix. In this example, water storage in the snow-
 1149 accounting store (solid line) starts (SE1) and ends (SE2) sooner than what would be expected from the SCA data
 1150 (dashed line).



Definition of the snow error (%):

$$\text{Crit4} = \text{SE} = \frac{1}{N_{\text{SCA}}} (\text{SE1} + \text{SE2})$$

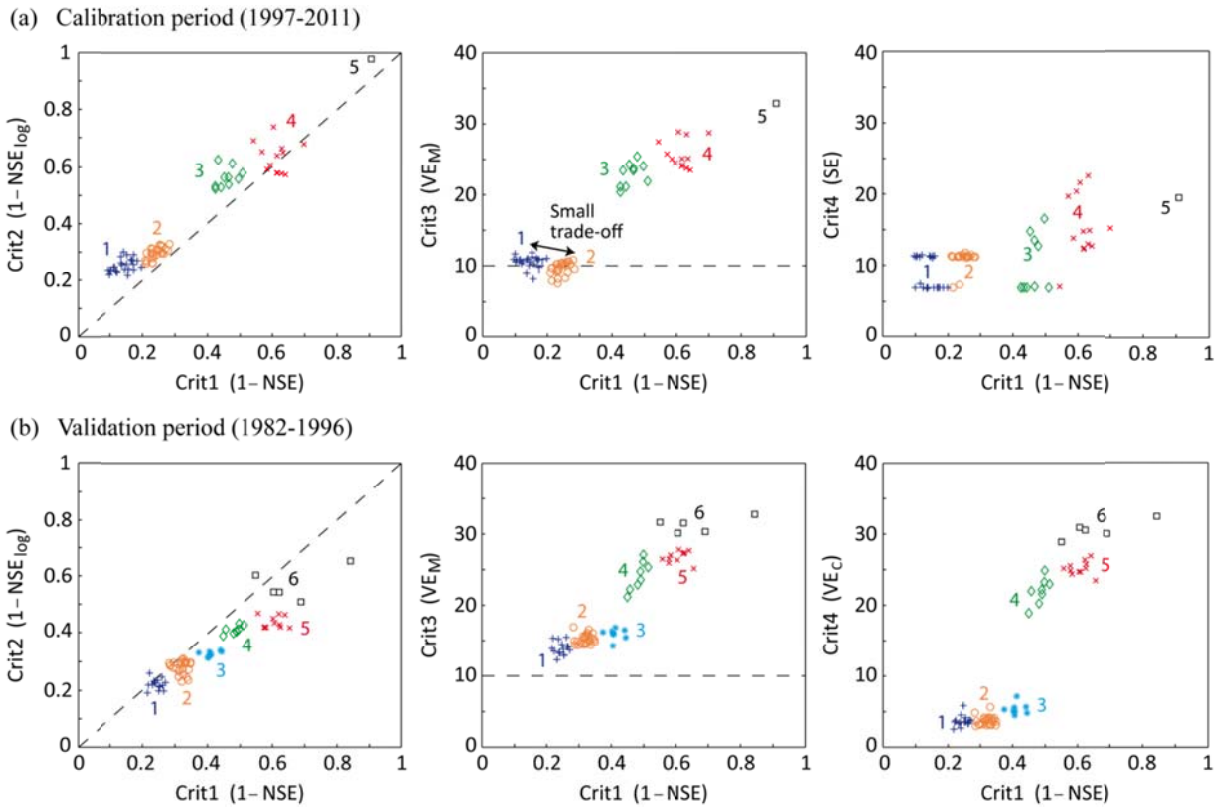
with N_{SCA} the number of days with available SCA observations

Confusion matrix (days) of the SE:

		SWE	
		> 0	= 0
SCA	> 0	No error	SE2
	= 0	SE1	No error

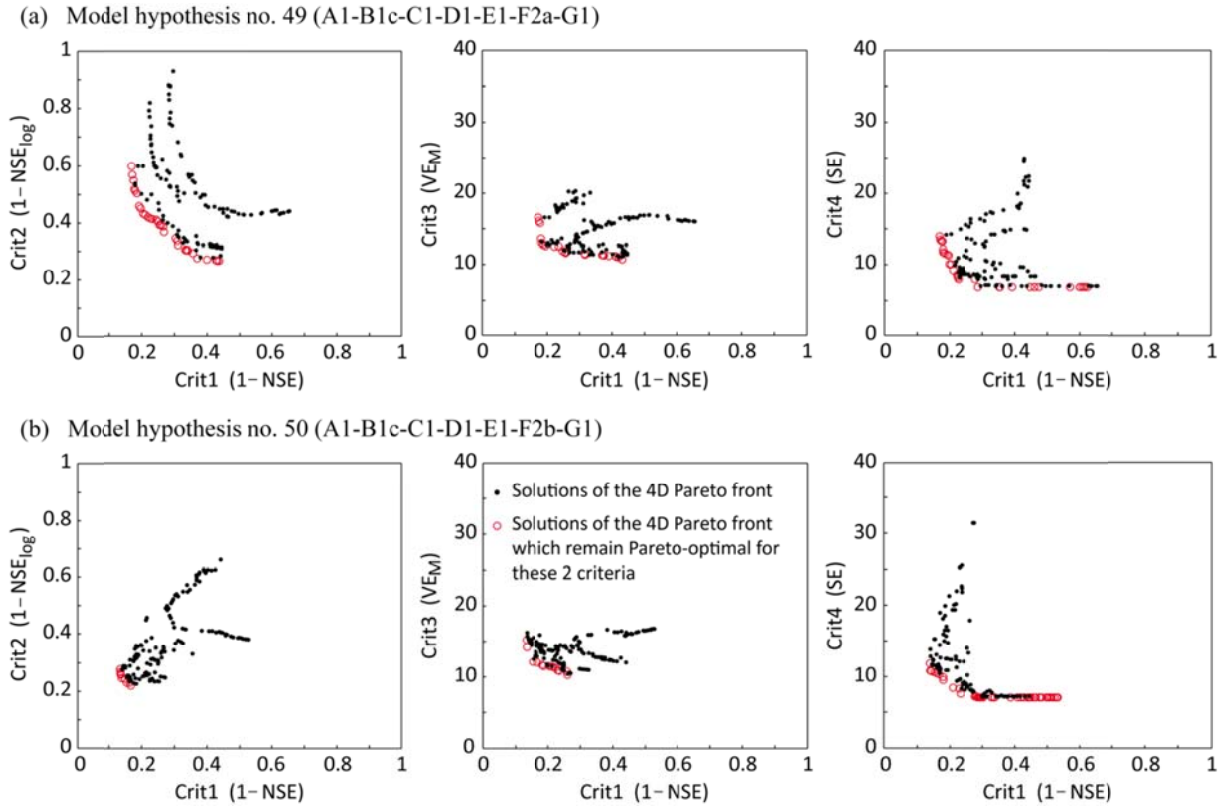
1151

1152 **Figure 5.** Projections of the clusters onto three possible planes of the objective space in calibration and
1153 validation. As explained in Sect 3.3., each point represents a different model hypothesis.
1154



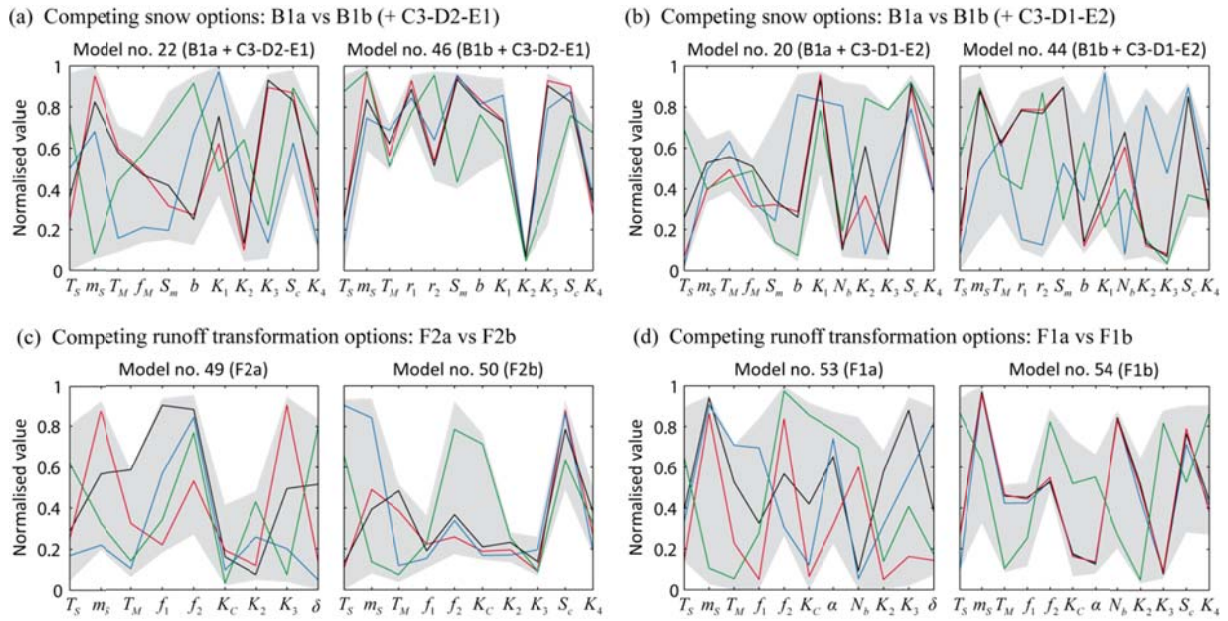
1155

1156 **Figure 6.** Projections of the Pareto fronts of model hypotheses (a) no. 49 (A1-B1c-C1-D1-E1-F2a) and (b) no.
 1157 50 (A1-B1c-C1-D1-E1-F2b) onto three possible two-dimensional subspaces of the objective space.



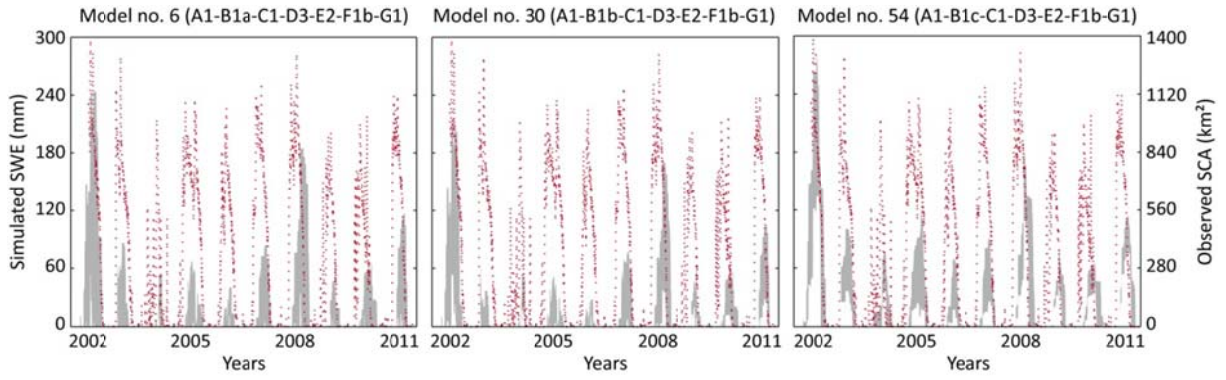
1158

1159 **Figure 7.** Estimated normalized ranges of the Pareto-optimal sets of eight alternative model structures differing
 1160 in at least one of their components. The colored lines stand for the best solutions obtained in calibration with
 1161 respect to the high flow criterion (in black), the low flow criterion (in red), the mean annual volume error (in
 1162 blue) and the snow error (in green).
 1163



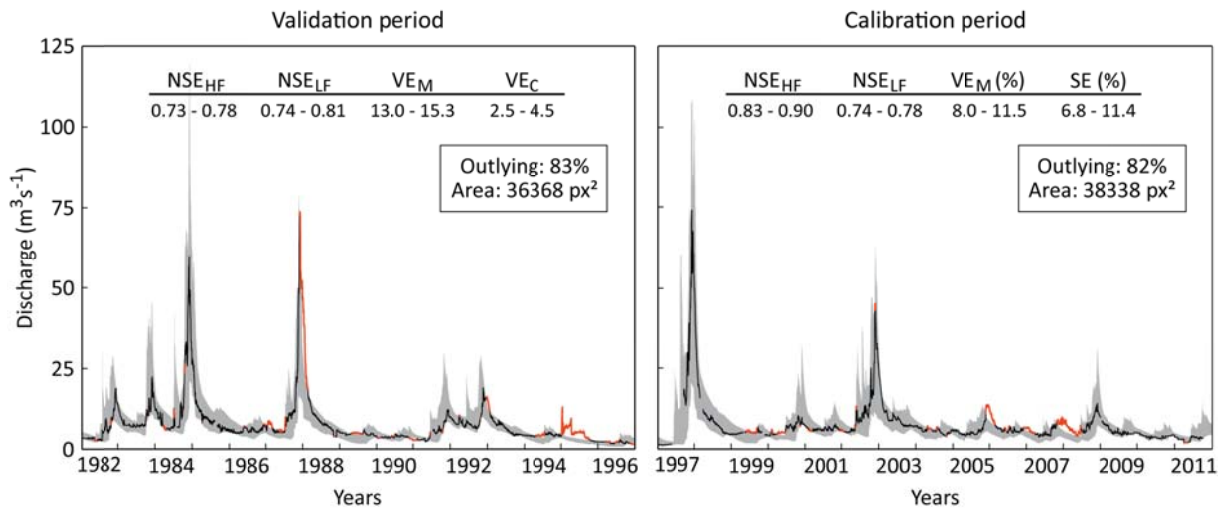
1164

1165 **Figure 8.** Comparison of MODIS-based SCA data (red dashed lines) with the SWE simulations (shaded areas)
1166 of models no. 6, 30 and 54. The shaded area corresponds to the range of SWE simulations obtained from the
1167 Pareto sets of these models.



1168

1169 **Figure 9.** Comparison of observed daily discharge at Rivadavia with the overall uncertainty envelope obtained
 1170 by combining the Pareto-envelopes of 8 model structures. These structures have been selected among the 14
 1171 members of Cluster 1 in both calibration and validation so as to minimize the uncertainty envelope area (Area, in
 1172 pixels²) while holding constant the number of outlying observations (Outlying, in %). The red parts indicate
 1173 potential errors in the model structures or observed data.



1174