

1 **Sensitivity analysis of prior model probabilities and the value of prior knowledge in the**
2 **assessment of conceptual model uncertainty in groundwater modelling**

3

4 Rodrigo Rojas^{a*}, Luc Feyen^b, Alain Dassargues^{a,c}

5

6

7 * Corresponding author

8 ^a Applied geology and mineralogy, Department of Earth and Environmental Sciences,
9 Katholieke Universiteit Leuven

10 Celestijnenlaan 200 E, B-3001 Heverlee, Belgium

11 Tel.: +32 016 326449; fax: +32 016 326401.

12 E-mail address: Rodrigo.RojasMujica@geo.kuleuven.be

13

14 ^b European Commission - DG Joint Research Centre (JRC), Institute for Environment and
15 Sustainability, Land management and natural hazards unit

16 TP261, Via E. Fermi 2749, 21027 Ispra (Va), Italy

17 Tel.: +39 0332 789258; fax: +39 0332 786653

18 E-mail address: luc.feyen@jrc.it

19

20 ^c Hydrogeology and environmental geology, Department of Architecture, Geology,
21 Environment, and Constructions (ArGEnCo), Université de Liège

22 B.52/3 Sart-Tilman, B-4000 Liège, Belgium

23 Tel.: +32 4 3662376; fax: +32 4 3669520

24 E-mail address: alain.dassargues@geo.kuleuven.be

1 **Abstract**

2 A key point in the application of multi-model Bayesian averaging techniques to assess the
3 predictive uncertainty in groundwater modelling applications is the definition of prior model
4 probabilities, which reflect the prior perception about the plausibility of alternative models. In
5 this article we analyze the influence of prior knowledge and prior model probabilities on
6 posterior model probabilities, multi-model predictions and conceptual model uncertainty
7 estimations. The sensitivity to prior model probabilities is assessed using an extensive
8 numerical analysis in which the prior probability space of a set of plausible
9 conceptualizations is discretised to obtain a large ensemble of possible combinations of prior
10 model probabilities. Additionally, we assess the value of prior knowledge about alternative
11 models in reducing conceptual model uncertainty by considering three example knowledge
12 states, expressed as quantitative relations among the alternative models. A constrained
13 maximum entropy approach is used to find the set of prior model probabilities that
14 correspond to the different prior knowledge states. For illustrative purposes, we employ a 3-
15 dimensional hypothetical setup approximated by 7 alternative conceptual models. Results
16 show that posterior model probabilities, leading moments of the predictive distributions and
17 estimations of conceptual model uncertainty are very sensitive to prior model probabilities,
18 indicating the relevance of selecting proper prior probabilities. Additionally, including proper
19 prior knowledge improves the predictive performance of the multi-model approach,
20 expressed by reductions of the multi-model prediction variances up to 60%. However, the
21 ratio between-model to total variance does not substantially decrease. This suggests that the
22 contribution of conceptual model uncertainty to the total variance can not be further reduced
23 based only on prior knowledge about the plausibility of alternative models. These results
24 advocate including proper prior knowledge about alternative conceptualizations in
25 combination with extra conditioning data to further reduce conceptual model uncertainty in
26 groundwater modelling predictions.

1 **Keywords**

- 2 Multi-model prediction, uncertainty assessment, maximum entropy, prior knowledge,
- 3 conceptual model uncertainty

1 **1. Introduction and scope**

2 Groundwater modelling has become an essential part of groundwater management and
3 accurate model predictions are required to ensure an acceptable degree of confidence in
4 model results. However, incomplete knowledge about the geological setting and scarce or
5 prone to error information about model parameters, boundary conditions and input data,
6 render the predictions of groundwater dynamics and pollutant transport uncertain. Practice,
7 on the other hand, suggests that once a conceptual model is successfully calibrated its results
8 are rarely questioned and the conceptual model is assumed to be correct (Bredehoeft, 2005;
9 Hojberg and Refsgaard, 2005). However, a successful calibration does not guarantee the
10 correctness of the conceptual model. Rather, many parameter sets together with different
11 conceptual models may produce equally good results in a calibration process (Bredehoeft,
12 2003; Harrar *et al.*, 2003; Carrera *et al.*, 2005). In this sense, relying on a single hydrological
13 concept will likely produce biased and under-dispersive predictions due to neglecting
14 conceptual model uncertainty (Neuman, 2003).

15

16 In recent years, a number of multi-model methods have been proposed to address the problem
17 of conceptual model uncertainty in hydrological modelling (Neuman, 2003; Poeter and
18 Anderson, 2005, Ajami *et al.*, 2005; Refsgaard *et al.*, 2006). These methods seek to obtain
19 consensus predictions from a set of plausible models by linearly combining individual model
20 predictions. One such approach is Bayesian Model Averaging (BMA) (Draper, 1995; Hoeting
21 *et al.*, 1999), which weights the predictions of competing models by their corresponding
22 posterior model probability, representing each model's relative skill to reproduce system
23 behaviour in the training period. Hence, BMA weights are tied directly to individual model
24 performance. Several studies applying the method to a range of different problems have
25 demonstrated that BMA produces more accurate and reliable predictions than other existing
26 multi-model techniques (e.g., Raftery and Zheng, 2003; Ye *et al.*, 2004; Ajami *et al.*, 2005).

1 In the field of groundwater hydrology, applications of BMA have been rare. Neuman (2003)
2 proposed the Maximum Likelihood Bayesian Model Averaging (MLBMA) method, which is
3 an approximation of BMA that relies on maximum likelihood parameter estimation and
4 expanding around these values through Monte Carlo simulation. Ye *et al.*, (2004) expanded
5 upon the theoretical framework of MLBMA and applied it to model the log permeability in
6 unsaturated fractured tuff using alternative variogram models.

7

8 Rojas *et al.*, (2008) proposed a methodology to assess uncertainty in predictions of
9 groundwater models arising from errors in the model structure, forcing data and parameter
10 estimates by integrating the Generalized Likelihood Uncertainty Estimation (GLUE) (Beven
11 and Binley, 1992) methodology with BMA. The methodology is based on the concept that
12 there exist many good simulators of the system that may be located in different regions of the
13 combined model, input and parameter space, given the data at hand. For a set of plausible
14 system conceptualizations, input and parameter realizations are sampled from the joint prior
15 input and parameter space. A likelihood measure is then calculated for each simulator based
16 on its ability to reproduce system state variable observations. The integrated likelihood of
17 each conceptual model is obtained by integrating the likelihood of the different simulators
18 over the input and parameter space. The integrated likelihoods are consequently used in BMA
19 to weight the model predictions to obtain ensemble predictions. Key advantages of this
20 methodology are that: (i) there is no restriction on the diversity of conceptual models or on
21 the level of uncertainty in the input data or parameters that can be included; (ii) it does not
22 rely on a single optimum set of (calibrated) parameter values, hence, avoiding biased
23 parameter estimates that compensate for errors in model structure, input data and
24 measurement errors; (iii) it allows for different ways of expressing the likelihood of a
25 simulator (including a formal Bayesian one), hence allowing different types of knowledge to
26 be incorporated (quantitative as well as qualitative); and (iv) it is Bayesian in nature, which

1 provides a formal framework to incorporate prior knowledge about the model structures and
2 parameters, or to update the estimates should new information become available.

3

4 Rojas *et al.*, (2008) applied the methodology by considering 7 alternative conceptualizations
5 with increasing complexity to represent a 3-dimensional synthetic example consisting of 2
6 aquifers separated by an aquitard. An extensive numerical analysis showed that neglecting
7 conceptual model uncertainty results in biased and overly conservative predictions. However,
8 two important aspects concerning the application of the methodology remained unanswered;
9 first, the sensitivity of posterior model probabilities, multi-model groundwater predictions,
10 and conceptual model uncertainty estimations to different sets of prior model probabilities;
11 and, second, the value of prior knowledge about the alternative conceptualizations to further
12 reduce conceptual model uncertainty. We address these two points in this article.

13

14 In Bayesian inference two basic interpretations can be given to prior probability distributions.
15 First, in the *population* interpretation, a prior distribution represents a population of possible
16 parameter values from which a potential candidate is to be drawn. Second, in the more
17 subjective *state of knowledge* interpretation, the guiding principle is that knowledge (and
18 uncertainty) about a given parameter must be expressed as if the value of that parameter
19 could be thought of as a random realization from the prior probability distribution (Gelman *et*
20 *al.*, 2004, p. 39), i.e., prior probability distributions can be interpreted as a formal
21 representation of knowledge (uncertainty) about a given parameter. More importantly, there
22 is no unique prior probability distribution for representing this knowledge (uncertainty) (Kass
23 and Wasserman, 1996).

24

25 In Bayesian literature, different methods to assign prior probability distributions to different
26 classes of problems can be found. We do not wish to provide a complete overview of these
27 methods but refer the reader to Kass and Wasserman (1996) for an excellent review.

1 A key point when adopting a prior probability distribution is the influence of this distribution,
2 after updating, on the results. Two general courses of action can be mentioned to alleviate
3 this influence. First, with increasing data availability, prior probability distributions are
4 expected to have less influence on inferences about parameters and predicted variables (Kass
5 and Wasserman, 1996). Thus, one strategy consists in collecting as much data as possible to
6 overcome the influence of the prior probability distributions. For most groundwater
7 modelling applications, however, obtaining enough data to overrule the effects of prior model
8 probabilities may in many cases be cost prohibitive. Second, one can assign non-informative
9 prior probability distributions, with the uniform distribution being the most common case,
10 hoping that information contained in the data will dominate the form of the resulting posterior
11 distribution. Consequently, reported multi-model methodologies used in groundwater
12 modelling have employed, generally, a uniform prior model probability distribution reflecting
13 no prior preference on the plausibility of alternative conceptual models (see, e.g., Meyer *et*
14 *al.*, 2007). This is also the approach followed by Rojas *et al.*, (2008).

15

16 Panels of experts, prior elicitation, and theoretical or empirical grounds, on the other hand,
17 can be helpful in defining suitable prior model probabilities based on expert knowledge (see,
18 e.g., Ye *et al.*, 2006). These prior model probabilities are inherently subjective, i.e., they
19 reflect preference over a particular conceptualization and, probably, other group of experts
20 will arrive to different prior model probabilities based on different grounds. In this context,
21 we stick to the idea expressed by Ghosh *et al.*, (2006, p. 55) who stated that whenever prior
22 information is available, an attempt to use a prior probability distribution reflecting that prior
23 knowledge should be used as far as possible.

24

25 Given that there is no unique way to express the prior knowledge about alternative conceptual
26 models, due mainly to the subjective nature of the task, a procedure to select among potential
27 sets of prior model probabilities is required. Ye *et al.*, (2005) recently proposed an approach

1 aimed to find a set of prior model probabilities that maximizes Shannon’s entropy (Shannon,
2 1948) subject to a series of constraints. Hereby, the constraints reflect prior knowledge about
3 the alternative conceptualizations. The key idea behind this approach is that uncertainty
4 represents “potential information” in the sense that when a random variable takes on a value
5 we gain information and lose uncertainty (Applebaum, 1996). In this sense, Shannon’s
6 entropy measures the amount of information contained in the set of prior model probabilities.
7 Therefore, less informative sets will have a higher entropy compared with more informative
8 sets since a larger amount of information can be gained in the first. For example, when the set
9 of prior model probabilities corresponds to the uniform prior distribution, i.e., all alternative
10 conceptual models have equal prior probabilities, we are at a state of maximum uncertainty
11 and entropy is at its maximum. When a more informative set of prior model probabilities is
12 available the entropy will be lower.

13

14 In the case that several sets of constraints reflecting different prior knowledge about the
15 conceptual models are proposed, the problem translates into a min-max choice, i.e., to find
16 the set of prior model probabilities that maximizes Shannon’s entropy subject to the
17 respective constraints, but which is minimum among different proposed sets. Solving this
18 min-max problem, however, does not guarantee optimum predictive performance. To
19 overcome this problem, Ye *et al.*, (2005) propose to follow one of the following two
20 approaches: (1) when enough data are available to perform a meaningful model
21 (cross)validation, they advocate selecting a posteriori the set that outperforms based on
22 suitable performance criteria (see, e.g., Liang *et al.*, 2001 for examples on performance
23 criteria); (2) when there is not enough data available to estimate meaningful posterior
24 measures of model quality, they advocate selecting the min-max set that, additionally,
25 maximizes the likelihood for the ensemble of alternative conceptual models.

1 In this work, we conduct a numerical experiment to analyze the sensitivity of the posterior
2 model probabilities, the groundwater multi-model predictions, and the conceptual model
3 uncertainty estimations to prior model probabilities. To this end, the prior probability space
4 of the alternative conceptual models is discretised in equidistant intervals and all possible
5 combinations of prior model probabilities for the set of conceptualizations are formed, given
6 that the sum of the prior model probabilities for each combination equals 1.

7

8 Furthermore, we extend upon the work of Ye *et al.*, (2005) and assess the value of prior
9 knowledge about the plausibility of alternative conceptualizations in reducing conceptual
10 model uncertainty. To this end we employ the constrained maximum entropy approach
11 proposing (out of the ensemble of discrete sets of prior model probabilities) three sets of prior
12 model probabilities that reflect the following knowledge states: (i) a non-informative case
13 about the plausibility of alternative conceptualizations, i.e., alternative conceptual models
14 have equal prior probabilities; (ii) relevant and proper prior knowledge about the plausibility
15 of alternative conceptualizations, i.e., alternative conceptual models receive higher prior
16 probabilities as they approach a “true” 3-dimensional hypothetical setup; and (iii) improper
17 prior knowledge about the plausibility of alternative conceptualizations, i.e., alternative
18 conceptual models receive prior probabilities that are inconsistent as they approach the “true”
19 3-dimensional hypothetical setup. Results obtained using the three optimized sets of prior
20 model probabilities are compared to find the set that outperforms in terms of predictive
21 capacity and to assess the value of this prior knowledge to further reduce conceptual model
22 uncertainty.

23

24 The remainder of this paper is organized as follows. In section 2, we provide a condensed
25 overview of the combined GLUE-BMA methodology. Section 3 details a 3-dimensional
26 hypothetical aquifer system that is used to illustrate the methodology and to assess the
27 sensitivity of the groundwater multi-model predictions. Implementation details are described

1 in section 4. In this section, we elaborate on the different conceptual models, input and
2 parameter uncertainty, the methodology to account for the sensitivity of the results due to
3 different discrete sets of prior model probabilities and the constrained maximum entropy
4 method to assess suitable sets of prior model probabilities in agreement with prior
5 knowledge. Results are discussed in section 5 and a summary of conclusions is presented in
6 section 6.

7

8 **2. Materials and Methods**

9 To render the article self-contained sections 2.1 and 2.2 elaborate on the basis of GLUE and
10 BMA methodologies, respectively. For a detailed description the reader is referred to Rojas *et*
11 *al.*, (2008).

12

13 **2.1. Generalized Likelihood Uncertainty Estimation (GLUE) methodology**

14 GLUE is a Monte Carlo simulation technique based on the concept of equifinality (Beven and
15 Freer, 2001). It rejects the idea of a single correct representation of a system in favour of
16 many acceptable system representations (Beven, 2005). For each potential system simulator,
17 sampled from a prior set of possible system representations, a likelihood measure is
18 calculated which reflects its ability to simulate the system responses, given the available
19 training data \mathbf{D} . Simulators that perform below a subjectively defined rejection criterion are
20 discarded from the further analysis and likelihood measures of retained simulators are
21 rescaled so as to render the cumulative likelihood equal to 1. Ensemble predictions are based
22 on the predictions of the retained set of simulators, weighted by their respective rescaled
23 likelihood.

24

25 Likelihood measures used in GLUE must be seen in a much wider sense than the formal
26 likelihood functions used in traditional statistical estimation theory (Binley and Beven, 2003).
27 These likelihoods are a measure of the ability of a simulator to reproduce a given set of

1 training data, therefore, they represent an expression of belief in the predictions of that
2 particular simulator rather than a formal definition of probability. However, GLUE is fully
3 coherent with a formal Bayesian approach when the use of a classical likelihood function is
4 justifiable (see, e.g., Romanowicz *et al.*, 1994).

5

6 In the work of Rojas *et al.*, (2008) no significant differences were observed in the estimation
7 of posterior model probabilities, predictive capacity and conceptual model uncertainty when
8 using a Gaussian, a model efficiency or a Fuzzy-type likelihood function. The analysis in this
9 work is therefore confined to a traditional Gaussian likelihood function $L(\mathbf{M}_k, \boldsymbol{\theta}_l, \mathbf{Y}_m | \mathbf{D})$,

10 where M_k is the k -th conceptual model (or model structure) included in the finite and
11 discrete ensemble of alternative conceptualizations \mathbf{M} , $\boldsymbol{\theta}_l$ is the l -th parameter vector, \mathbf{Y}_m is
12 the m -th input data vector, and \mathbf{D} is the observed system variable vector.

13

14 **2.2. Bayesian Model Averaging (BMA)**

15 BMA provides a coherent framework for combining predictions from multiple conceptual
16 models to attain a more realistic and reliable description of the total prediction uncertainty. It
17 yields consensus predictions by weighing predictions from competing models based on their
18 relative skill, with predictions from better performing models receiving higher weights than
19 those of worse performing models. BMA avoids having to choose one model over the others,
20 instead, observed data \mathbf{D} give the competing models different weights (Wasserman, 2000).

21

22 Following the notation of Hoeting *et al.*, (1999), if Δ is a quantity to be predicted, the BMA
23 predictive distribution of Δ is given by

24

$$25 \quad p(\Delta | \mathbf{D}) = \sum_{k=1}^K p(\Delta | \mathbf{D}, M_k) p(M_k | \mathbf{D}). \quad (1)$$

1 Equation 1 is an average of the posterior distributions of Δ under each alternative conceptual
 2 model considered, $p(\Delta | \mathbf{D}, M_k)$, weighted by their posterior model probability, $p(M_k | \mathbf{D})$.
 3 This latter term reflects how well model k fits the observed data \mathbf{D} and can be computed
 4 using Bayes' rule

$$6 \quad p(M_k | \mathbf{D}) = \frac{p(\mathbf{D} | M_k) p(M_k)}{\sum_{k'=1}^K p(\mathbf{D} | M_{k'}) p(M_{k'})} \quad (2)$$

7
 8 where $p(M_k)$ is the prior probability of model k , and $p(\mathbf{D} | M_k)$ is the integrated likelihood
 9 of the model k .

10
 11 The leading moments of the BMA prediction of Δ are given by Draper (1995)

$$13 \quad E[\Delta | \mathbf{D}] = \sum_{k=1}^K E[\Delta | \mathbf{D}, M_k] p(M_k | \mathbf{D}) \quad (3)$$

$$15 \quad \begin{aligned} Var[\Delta | \mathbf{D}] = & \sum_{k=1}^K Var[\Delta | \mathbf{D}, M_k] p(M_k | \mathbf{D}) \\ & + \sum_{k=1}^K (E[\Delta | \mathbf{D}, M_k] - E[\Delta | \mathbf{D}])^2 p(M_k | \mathbf{D}) \end{aligned} \quad (4)$$

16
 17 From equations 1 and 2 it is seen that estimations of posterior model probabilities (weights)
 18 and, subsequently, estimations of the first two leading moments of the BMA predictive
 19 distribution (equations 3 and 4), are functions of the prior model probabilities assigned to the
 20 alternative conceptual models. From equation 4 it is seen that the variance of the BMA
 21 predictions consists of two terms; the first representing the within-model variance and the

1 second representing the between-model variance (variance due to conceptual model
2 uncertainty).

3

4 **2.3. Combining GLUE and BMA**

5 Combining GLUE and BMA involves the following sequence of steps

- 6 1. Based on prior and expert knowledge about the site, a suite of alternative conceptual
7 models is proposed.
- 8 2. Realistic prior ranges are defined for the input and parameter vectors under each plausible
9 model structure.
- 10 3. A likelihood measure and rejection criteria are defined.
- 11 4. For the suite of alternative conceptual models, input and parameter values are sampled
12 from the prior ranges to generate possible simulators of the system.
- 13 5. A likelihood measure is calculated for each simulator based on the agreement between the
14 simulated and observed system response.
- 15 6. Simulators that are not in agreement with the selected rejection criterion are discarded
16 from the analysis by setting their likelihood to zero.
- 17 7. For each conceptual model M_k , a subset A_k of simulators with likelihood
18 $p(\mathbf{D}|M_k, \boldsymbol{\theta}_l, \mathbf{Y}_m) = L(M_k, \boldsymbol{\theta}_l, \mathbf{Y}_m | \mathbf{D})$ is retained. Steps 4-6 are repeated until the
19 hyperspace of possible simulators is adequately sampled, i.e., when the conditional
20 distributions of predicted state variables based on the likelihood weighted simulators in
21 the subset A_k converge to a stable distribution for each of the conceptual models M_k .
- 22 8. The integrated likelihood of each conceptual model M_k is approximated by summing the
23 likelihood weights of the retained simulators in subset A_k , or

24

$$25 \quad p(\mathbf{D}|M_k) \approx \sum_{l,m \in A_k} L(M_k, \boldsymbol{\theta}_l, \mathbf{Y}_m | \mathbf{D}) \quad (5)$$

1 9. The posterior model probabilities are then obtained by normalizing the integrated model
 2 likelihoods such that they sum up to 1,

3

$$4 \quad p(M_k | \mathbf{D}) \approx \frac{\sum_{A_k} L(M_k, \boldsymbol{\theta}_l, \mathbf{Y}_m | \mathbf{D}) p(M_k)}{\sum_{j=1}^K \sum_{l, m \in A_j} L(M_j, \boldsymbol{\theta}_l, \mathbf{Y}_m | \mathbf{D}) p(M_j)} \quad (6)$$

5

6 10. After normalization of the likelihood weighted predictions under each individual model
 7 (such that the cumulative likelihood under each model equals 1) a multi-model prediction
 8 is obtained with equation 1 using the weights obtained with equation 6.

9

10 Details about the implementation of the methodology, applied to the 3-dimensional
 11 hypothetical setup described in the next section, are presented in Section 4.

12

13 **3. Three-dimensional hypothetical setup**

14 For illustrative purposes, we employ a 3-dimensional hypothetical setup for which the true
 15 conditions are known (Figure 1). Lateral dimensions are 5000 m (E-W) by 3000 m (N-S)
 16 discretised in 25 m by 25 m grid cells. The system extents over 60 m in the vertical direction,
 17 with undisturbed layer thicknesses of 35 m (upper aquifer), 5 m (middle aquitard) and 20 m
 18 (lower aquifer). We assume statistically homogeneous deposits with a constant mean
 19 hydraulic conductivity K (see Table 1). Smaller-scale variability is represented using the
 20 theory of random space functions, adopting isotropic exponential covariance functions for log
 21 K in all layers. The spatial distribution of the hydraulic conductivity in the layers of the
 22 example setup, as well as any other realization of the hydraulic conductivity field used in this
 23 work, is generated using the sequential Gaussian simulation (sGsim) algorithm of the
 24 Geostatistical Software Library (Deutsch and Journel, 1998). Parameters of the covariance
 25 function of log K for the different layers are presented in Table 1.

1 Simulation of steady-state flow is performed using Modflow-2000 (Harbaugh *et al.*, 2000).
2 At the north and south boundaries, as well as at the bottom of the lower layer, zero gradient
3 conditions are imposed. A uniform recharge of $1.4 \times 10^{-4} \text{ m d}^{-1}$ is applied to the top layer. At
4 the west boundary a constant head $h = 46 \text{ m}$ is defined. The east side of the domain is
5 bounded by a 10 m-wide river with a constant stage of 25 m. The river bottom is at 20 m,
6 defining a constant river water depth of 5 m. It is underlain by 5 m-thick sediments with a
7 vertical hydraulic conductivity of 0.1 m d^{-1} . Five pumping wells are distributed in the area
8 producing a total of $2450 \text{ m}^3 \text{ d}^{-1}$ from the lower aquifer (Figure 1). An evapotranspiration
9 zone, delineated by the polygon in Figure 1, is defined with an evapotranspiration surface
10 elevation at 43 m, an evapotranspiration rate of $1.37 \times 10^{-3} \text{ m d}^{-1}$ and an extinction depth of 5
11 m.

12
13 The resulting “true” groundwater head distribution for the top layer is presented as an overlay
14 in Figure 1. The ambient background gradient from west to east is strongly influenced by the
15 drawdown around pumping wells, the evapotranspiration zone as well as by local effects of
16 spatially varying hydraulic conductivity. From the “true” groundwater head distribution for
17 layer 1, values are selected at the 16 locations defined by the observation wells in Figure 1,
18 which are used to estimate the likelihood weights in the evaluation of different simulators.

19

20 **4. Implementation of the methodology and numerical analysis**

21 **4.1. Implementation of the GLUE-BMA approach**

22 We consider 7 alternative conceptual models with increasing complexity to describe the 3-
23 dimensional hypothetical setup described in section 3, namely: (1), (2) and (3) one-layer
24 models with mean K and spatial correlation law of layer 1 (1Lhtg-L1), layer 2 (1Lhtg-L2) and
25 layer 3 (1Lhtg-L3) of the hypothetical setup, respectively; (4) a one-layer model with average
26 mean K and spatial correlation (1Lhtg-AVG); (5) a two-layer model with mean K and spatial
27 correlation taken from layer 1 and layer 3 (2Lhtg); (6) a two-layer quasi-three dimensional

1 model with mean K and spatial correlation taken from layer 1 and layer 3, and mean K of
2 layer 2 used to define the aquitard (2LQ3Dhtg); and (7) a three-layer model based on the
3 spatial K distributions of layer 1, layer 2 and layer 3 (3Lhtg). All conceptual models comprise
4 a total aquifer thickness of 60 m and are forced by identical types of boundary conditions.
5

6 The dimensionality of the analysis is confined by considering uncertainty only in the input
7 variables and parameters related to the evapotranspiration process, lateral boundary
8 conditions, river description and recharge process, i.e., input variables and parameters that are
9 common to all setups. Values are sampled from uniform prior distributions for the unknown
10 inputs and parameters with ranges defined in Table 2. Unconditional realizations of the
11 hydraulic conductivity field are generated with the same mean K and spatial correlation law
12 as the respective layers in the hypothetical setup (Table 1). For the 1Lhtg-AVG
13 conceptualization the average of these values is used.
14

15 For the simulation, parameter and input vectors sampled using a Latin Hypercube Sampling
16 (LHS) scheme, are combined with unconditional hydraulic conductivity realizations and
17 consequently evaluated under each conceptual model. Based on the evaluation of a set of
18 initial runs, a rejection threshold is defined corresponding to a maximum allowable deviation
19 of 5 m at any of the 16 observation wells depicted in Figure 1. A point rejection threshold
20 rather than a global rejection threshold is chosen because under the latter criteria strong
21 deviations at certain locations (typically in the vicinity of pumping wells) may be offset by
22 small deviations at other wells. For each conceptual model, predictive distributions for the
23 sixteen observation wells depicted in Figure 1 and different components of the groundwater
24 budget (recharge inflows, groundwater inflows/outflows from the west boundary condition
25 (WBC), river gains, and evapotranspiration (EVT) outflows) are obtained from the ensemble
26 of likelihood weighted predictions. Sampling from the prior input and parameter space
27 continued until the first and second moment of these predictive distributions stabilized.

1 **4.2. Approach to assess sensitivity to prior model probabilities**

2 To analyze the sensitivity to different values of prior model probabilities, the prior model
3 probability space of the 7 alternative conceptualizations is discretised into 25 equidistant
4 intervals of 4% probability each. To avoid extremely low model probabilities that *reject with*
5 *high certainty* one of the proposed alternative conceptual models, the lowest probability
6 intervals are discarded from the analysis (this implies that the highest probability of a model
7 is $1-6*0.04 = 0.76$). From the remaining 19 probability intervals the lowest value of each
8 interval is retained, resulting in the following set of potential prior model probabilities for
9 each of the 7 alternative conceptual models $P = [0.04, 0.08, \dots, 0.76]$. Subsequently, all
10 combinations that fill the prior probability space conditional on \mathbf{M} , i.e., for which

11 $\sum_{k=1}^K p(\mathbf{M}_k) = 1$ (243 vectors of 7 elements), are formed. This yields a total of 132,861

12 potential discrete sets of prior model probabilities that are used to numerically analyze the
13 sensitivity of the posterior model probabilities (weights in equation 1), multi-model
14 predictions and conceptual model uncertainty estimation to prior model probabilities.

15

16 **4.3. Constrained maximum entropy approach to assess value of prior knowledge**

17 The value of prior knowledge about the plausibility of the 7 alternative conceptual models in
18 assessing conceptual model uncertainty is evaluated following a constrained maximum
19 entropy method (Ye *et al.*, 2005). The method aims to find discrete sets of prior model
20 probabilities that maximizes Shannon's entropy H (Shannon, 1948) given by

21

$$22 \quad H = - \sum_{k=1}^K p(\mathbf{M}_k) \log p(\mathbf{M}_k) \quad (7)$$

23

24 and subject to

$$\begin{aligned}
1 \quad & h_i = 0 && i = 1, \dots, I \\
& g_j = 0 && j = 1, \dots, J
\end{aligned} \tag{8}$$

2

3 where h_i and g_j represent quantitative relations that reflect prior knowledge about the
4 plausibility of the alternative conceptual models. In equation (7) $p(M_k)$ is the prior model
5 probability of the k -th conceptual model contained in the ensemble \mathbf{M} of dimension K . In the
6 case that alternative sets of constraints (reflecting different knowledge states) are proposed
7 and when not enough data are available to assess the quality of model results, Ye *et al.*,
8 (2005) advocate selecting the set that: (i) maximizes entropy H , (ii) presents a minimum
9 entropy among proposed sets and, (iii) maximizes the likelihood for \mathbf{M} given by the
10 normalizing term in equation 2.

11

12 For illustrative purposes we define 3 different prior knowledge states: (i) Prior Set 1,
13 corresponding to a set of uniform prior model probabilities $p(M_k) = 1/K$, reflecting a state
14 of complete ignorance about the plausibility of the alternative conceptual models; (ii) Prior
15 Set 2, corresponding to a set where alternative conceptual models receive higher prior
16 probability as they approach the 3-dimensional hypothetical setup described in section 3 and,
17 thus, reflecting relevant and proper prior knowledge about the alternative conceptualizations;
18 and (iii) Prior Set 3, corresponding to a set where prior model probabilities are inconsistent
19 with the degree of similarity between the alternative conceptual models and the 3-
20 dimensional hypothetical setup and, thus, reflecting improper prior knowledge about the
21 alternative conceptualizations.

22

23 We adopted the following set of constraints to reflect the information contained in the three
24 proposed prior knowledge states

$$\max_{p(M_k)} H = - \sum_{k=1}^K p(M_k) \log p(M_k)$$

1 Set 1:

$$h_1 = \sum_{k=1}^K p(M_k) - 1$$

$$g_1 \dots g_6 : p(M_1) = p(M_2) = p(M_3) = p(M_4) = p(M_5) = p(M_6) = p(M_7)$$

2

$$\max_{p(M_k)} H = - \sum_{k=1}^K p(M_k) \log p(M_k)$$

$$h_1 = \sum_{k=1}^K p(M_k) - 1$$

3 Set 2:

$$g_1 : p(M_1) - p(M_2) = 0$$

$$g_2 : p(M_2) - p(M_3) = 0$$

$$g_3 : p(M_4) - 2.0p(M_3) \geq 0$$

$$g_4 : p(M_5) - 1.5p(M_4) \geq 0$$

$$g_5 : p(M_6) - 2.5p(M_5) \geq 0$$

$$g_6 : p(M_7) - 1.1p(M_6) \geq 0$$

4

$$\max_{p(M_k)} H = - \sum_{k=1}^K p(M_k) \log p(M_k)$$

$$h_1 = \sum_{k=1}^K p(M_k) - 1$$

5 Set 3:

$$g_1 : p(M_1) - p(M_2) = 0$$

$$g_2 : p(M_2) - p(M_3) = 0$$

$$g_3 : 0.5p(M_3) - p(M_4) \geq 0$$

$$g_4 : 0.8p(M_4) - p(M_5) \geq 0$$

$$g_5 : 0.5p(M_5) - p(M_6) \geq 0$$

$$g_6 : 0.5p(M_6) - p(M_7) \geq 0$$

6

7 For Prior Set 1 (uniform prior model distribution) the solution to the optimization problem is

8 known to be $H = \log K = 1.95$ (see, e.g., Applebaum, 1996, p. 100) with $p(M_k) = 1/7$. For

9 Prior Set 2 and 3 the nonlinear optimization problem is solved numerically using a sequential

10 equality constrained quadratic programming method implemented in an R interface (Tamura,

11 2007) for the code DONLP2 (Spellucci, 1998). The result of these optimization problems are

1 three optimized sets of prior model probabilities for the 7 alternative conceptual models that
2 are in agreement with the quantitative relations (constraints) expressing the prior knowledge
3 states. The optimized values are presented in Table 3. These three sets of prior model
4 probabilities are samples from the full range of possible prior probability combinations,
5 approximated here by the ensemble of discrete sets. It is important to note that the values of
6 the constants in the constraints for Prior Set 2 and 3 were set as an example. Other values for
7 these constants would result in different prior model probabilities, however, still reflecting
8 prior knowledge. Consequently, the present analysis is conditional on the proposed ensemble
9 of alternative conceptual models, \mathbf{M} , and to the potential quantitative relations among them,
10 i.e., h_i and g_i .

11

12 **5. Results and discussion**

13 In the numerical analysis, for the alternative conceptual models 1Lhtg-L1 and 1Lhtg-L2 none
14 of the simulations were accepted, as all of them failed to meet the criterion of a maximum
15 allowable departure of 5 m from the observed heads. This suggests that approximating the
16 “true” 3-dimensional hypothetical setup using only information from layers 1 and 2 (see
17 Table 1) is not supported by the training data \mathbf{D} (i.e., observed head at 16 observation wells).
18 Hence, the posterior probability of these conceptual models was set to zero and they were
19 discarded from the posterior analysis.

20

21 **5.1. Sensitivity of posterior model probabilities to prior model probabilities**

22 The sensitivity of the posterior model probabilities to prior model probabilities for the 5
23 retained conceptual models is presented in Figure 2. In this figure, vertical columns represent
24 posterior model probabilities (estimated using equation 6) corresponding to the 132,861
25 nonzero discrete sets of prior model probabilities described in section 4.2. It can be seen that
26 the posterior model probabilities are sensitive to values of prior model probabilities for all the
27 retained models. It should be noted that the increase of the posterior model probabilities for

1 the 5 retained conceptual models, i.e., nearly all points lie above the bisector curve, is caused
2 by the fact that 2 out of 7 alternative conceptual models were discarded from the posterior
3 analysis based on the information contained in the training data \mathbf{D} . As a consequence, the
4 share in the prior probability space of the discarded conceptualizations is redistributed over
5 the 5 retained conceptual models when filling the posterior probability space (i.e., sum of
6 posterior probabilities should equal to 1). This explains why in most cases the posterior
7 probability is larger than the prior probability for the retained models. Notwithstanding, for
8 alternative conceptualizations 1Lhtg-L3 (Figure 2a) and 1Lhtg-AVG (Figure 2b) values of
9 posterior model probabilities below the bisector curve can be found, suggesting that less
10 weight is assigned a posteriori to these models. For alternative conceptual models 2Lhtg
11 (Figure 2c), 2LQ3Dhtg (Figure 2d) and 3Lhtg (Figure 2e), on the other hand, posterior model
12 probabilities are always higher than prior model probabilities, this being more noticeable for
13 model 3Lhtg.

14
15 From Figure 2 it is also seen that the uncertainty in the estimation of posterior model
16 probabilities (expressed by the range of the vertical columns) is maximum when there is no
17 clear preference a priori for a given conceptual model. On the contrary, the range of potential
18 values for posterior model probabilities is reduced when an alternative conceptual model is
19 preferred over the others.

20
21 Results for the three example sets of optimised prior model probabilities are also included in
22 Figure 2 and are summarized in Table 3. Results confirm that posterior model probabilities,
23 $p(M_k|\mathbf{D})$ are largely influenced by the selection of a set of prior model probabilities. For Prior
24 Sets 1 and 2, all retained models received more weight after conditioning. For Prior Set 2, on
25 the other hand, the posterior probability of the two retained one-layer models was smaller
26 than their respective prior probability, whereas the other 3 retained models received more
27 weight after conditioning. However, for all 3 sets, the relative increase of the posterior

1 probability compared to the prior probability is larger for the models approaching the true
2 setup.

3

4 **5.2. Sensitivity of the prior entropy, likelihood ratio and posterior entropy to prior** 5 **model probabilities**

6 The sensitivity of the prior entropy, likelihood ratio (with respect to the non-informative case)
7 and posterior entropy (calculated using equation 7 with $p(M_k|\mathbf{D})$ instead of $p(M_k)$) is
8 presented in Figure 3 for model 3Lhtg. It is seen in this figure that prior and posterior entropy
9 decreased when prior model probabilities of model 3Lhtg increased. Moreover, the likelihood
10 ratio (with respect to the non-informative case) tends to be maximized (Figure 3b) for a
11 maximum probability of model 3Lhtg. Consider, for example, a prior model probability of
12 0.76 for model 3Lhtg and, consequently, 0.04 for the 6 remaining models. Clearly, this set of
13 prior model probabilities is optimum (globally) in the sense that it minimizes posterior
14 entropy and it maximizes the likelihood ratio.

15

16 For the 3 example sets, the smallest maximum prior entropy, the smallest posterior entropy,
17 which can be interpreted as a measure of residual uncertainty after conditioning on the
18 training data \mathbf{D} (Ye *et al.*, 2005), and the largest likelihood ratio (1.34 times that of Prior Set
19 1) are obtained for Prior Set 2. On the contrary, the lowest likelihood ratio is observed for
20 Prior Set 3, which suggests that this set is not in agreement with the information contained in
21 the data and that it constitutes an improper expression of prior knowledge about the
22 alternative conceptual models. Hence, for the problem at hand, a reasonable choice for a
23 discrete set of prior model probabilities is to assign increasing probabilities in function of
24 proximity to the 3-dimensional hypothetical setup, i.e., Prior Set 2.

25

26 **5.3. Sensitivity of multi-model predictions and conceptual model uncertainty** 27 **estimations**

1 The sensitivity of the leading moments (estimated using equations 3 and 4) for model output
2 river gains and for three alternative conceptual models (1Lhtg-L3, 2Lhtg and 3Lhtg) is
3 presented in Figure 4. This figure shows that the posterior moments (plates a-f) of the
4 predictive distribution for river gains are rather sensitivity to prior model probabilities. It is
5 also seen that uncertainty in the estimation of the leading moments (expressed as the range of
6 the vertical columns) increased when the corresponding prior model probabilities decreased.
7 Additionally, when prior model probabilities for each alternative model increased, the leading
8 moments converged to different values. The latter suggests that when a model is preferred
9 over the others, i.e., relying only on a single conceptual model, predictions and uncertainty
10 estimations tend to be biased. Moreover, estimation of the leading moments tends to be
11 markedly more biased when prior model probabilities of simpler model 1Lhtg-L3 increased.
12
13 Plates g, h and i of Figure 4 show between-model variances for models 1Lhtg-L3, 2Lhtg and
14 3Lhtg, respectively, which are an expression of the conceptual model uncertainty. In general,
15 the contribution of conceptual model uncertainty to the total spread is sensitive to prior model
16 probabilities. Uncertainty in the estimation of between-model variance (expressed as the
17 range of the vertical columns) increased when prior model probabilities decreased. Moreover,
18 for the alternative conceptual models, between-model variances converged to different values
19 when corresponding prior model probabilities increased. It should be noted, however, that for
20 models 2Lhtg and 3Lhtg the converged values of between-model variances (2.1×10^3 and 2.8
21 $\times 10^3 \text{ [m}^3 \text{ d}^{-1}]^2$, respectively) were rather similar for a maximum prior model probability of
22 0.76. However, the ratio between-model to total variance was somewhat different (7% and
23 18%, for 2Lhtg and 3Lhtg, respectively) due to the difference in the estimation of total
24 variance for these models.
25
26 Figure 5 shows contour plots of the total variance and between-model variance (expressed as
27 a percentage of the total variance) for model outputs west boundary condition (WBC)

1 inflows, river gains and EVT outflows in the prior model probability space of 1Lhtg-L3
2 (simpler model) and 3Lhtg (model closer to the 3-dimensional hypothetical setup) when the 3
3 remaining alternative conceptual models approach a value near the uniform case (0.16). As
4 consequence, only 52% of the prior model probability space is left to be distributed in the
5 plates of Figure 5. More important than the actual values of the contour lines (which are
6 approximations since the true uniform case has a value of 0.143) is the shape of the surface
7 defined in the prior model probability space.

8
9 Plates a, b and c of Figure 5 show that the rate of change of the total variance (a measure of
10 sensitivity) is much larger in the prior space of model 1Lhtg-L3 (x-axis) compared to the
11 prior space of model 3Lhtg (y-axis). Hence, a more important reduction of the total variance
12 would be expected when prior model probabilities of 1Lhtg-L3 decrease. This suggests that,
13 for the problem at hand, to obtain more accurate multi-model predictions, simpler models
14 should receive less prior weight compared to more elaborated models. In addition, it is seen
15 from plates d, e and f that between-model variances does not fall below 5%, 20% and 12% of
16 the total variance and, on the other hand, they can reach values as large as 12%, 30% and
17 18% of the total variance for WBC inflows, river gains and EVT outflows, respectively.
18 Furthermore, the maximum contribution of between-model variances to total variances tends
19 to be located around the middle area of the figures, which is contrasting with the fact that the
20 non-informative case (uniform prior model probabilities) is not located in this area.

21
22 Overall, Figure 5 suggests that when a conceptual model tends to be preferred over the
23 others, between-model variance tends to be minimum. This is in agreement with previous
24 statements about under-dispersive properties of uncertainty estimations based on a single
25 model. On the contrary, between-model variance tends to be maximum when there is no clear
26 preference for a given conceptual model, suggesting that uncertainty estimations based on a
27 suite of alternative models are more spread. This seems logic since including alternative

1 conceptual models provides a more conservative assessment of uncertainty due to including
2 conceptual model uncertainty.
3
4 Figures 4 and 5 also include values for the three optimized sets of prior model probabilities.
5 Although posterior moments converged to different values for different conceptual models in
6 Figure 4, convergence was in agreement with the values obtained using Prior Set 2 when
7 models approached the “true” 3-dimensional hypothetical setup (see, e.g., plates c, f and i).
8 This supports the idea stated before that Prior Set 2 is a suitable choice to assign prior model
9 probabilities. This is also supported by the evidence provided by the data, which gave slightly
10 higher integrated model likelihood values to model 3Lhtg. It is also seen from Figures 4 and
11 5 that the posterior variance, with respect to the non-informative case (Prior Set 1),
12 significantly decreased when proper prior knowledge (Prior Set 2) was included in the
13 analysis. On the contrary, in the case of improper prior knowledge (Prior Set 3) a significant
14 increase of the total variance was observed. More importantly, between-model variances
15 (plates g, h and i of Figure 4) significantly decreased with respect to the non-informative case
16 (Prior Set 1) when proper prior knowledge (Prior Set 2) was included in the analysis,
17 indicating the value of prior knowledge in reducing conceptual model uncertainty.
18
19 Similar results were found for the other groundwater budget terms (Table 4). With respect to
20 Prior Set 1, total variances decreased between 40 and 60 % when the more informative Prior
21 Set 2 was used. On the contrary, total variances increased between 32 and 60% when
22 improper prior knowledge was included (Prior Set 3). Between-model variances decreased
23 for the informative Prior Set 2 by 50 up to 62% with respect to Prior Set 1. However, the
24 relative contribution of between-model variance to the total variance did not substantially
25 decrease. For example, for EVT outflows obtained using Prior Set 1, the contribution of
26 between-model to total variance is 0.15 whereas for Prior Set 2 this ratio is 0.14. The largest
27 reduction in the contribution of between-model to total variance for Prior Set 2 is observed

1 for river gains; from 0.26 to 0.2. This suggests that the contribution of conceptual model
2 uncertainty to total uncertainty can not be further reduced based only on prior knowledge
3 about the plausibility of alternative conceptualizations. This indicates that other sources of
4 information or conditioning data should be included to further reduce this component of total
5 variance.

6

7 For Prior Set 3 the between-model variances for WBC inflows and river gains increased,
8 whereas for recharge inflows, WBC outflows and EVT outflows, between-model variances
9 decreased compared to Prior Set 1. This erratic behaviour in the between-model variances
10 estimated using Prior Set 3 is explained by Figure 5.

11

12 **5.4. Value of prior knowledge about alternative conceptualizations in the goodness of** 13 **GLUE-BMA predictions**

14 Summary statistics of the posterior predictive distributions for the groundwater budget terms
15 as a function of the optimized sets of prior model probabilities are presented in Figure 6. In
16 this figure maximum values are truncated to enhance visual comparison. Observed values for
17 the groundwater budget terms, obtained from the 3-dimensional hypothetical setup, are
18 captured by the inter-quartile range of Prior Set 1 and Prior Set 2. On the contrary, observed
19 values for WBC outflows, river gains and EVT outflows are not captured by the inter-quartile
20 range of Prior Set 3. Comparing the optimized sets for each plate in Figure 6, Prior Set 2
21 outperforms the other sets since the median values are closer to the observed values and its
22 inter-quartile range is more concentrated, indicating less residual uncertainty after observing
23 data **D** and incorporating prior knowledge. Hence, this suggests that multi-model predictions
24 obtained using the GLUE-BMA approach in combination with proper prior knowledge (Prior
25 Set 2) outperforms multi-model predictions obtained using sets reflecting a non-informative
26 case (Prior Set 1) and improper prior knowledge (Prior Set 3).

1 GLUE-BMA predictions for groundwater heads at the locations depicted in Figure 1 are
2 presented in Figure 7. The predictive mean and standard deviation are estimated using
3 equations 3 and 4, respectively. The more pronounced differences in the mean predicted head
4 are observed for observation wells Obs-8, Obs-13, Obs-14, Obs-15 and Obs-16. It is
5 interesting to note that, for these observation wells, observed heads are captured by the
6 interval (± 1 standard deviation) defined around the predicted mean value using Prior Set 2.
7 On the contrary, observed heads are not captured by the interval defined using Prior Set 1 and
8 Prior Set 3. The exception to this is observation well Obs-2, in which none of the optimized
9 sets was able to capture the observed head. It is also shown in Figure 7 that for some
10 observation wells the standard deviations obtained using Prior Set 3 are slightly smaller
11 compared to those obtained with the other optimized sets. However, this gain in accuracy is
12 irrelevant since observed heads are not captured by the intervals defined using Prior Set 3 in 7
13 out of 16 observation wells. Therefore, an over-confident and biased prediction of the
14 observed heads is obtained when improper prior knowledge (Prior Set 3) is used.

15

16 These results confirm that, for the problem at hand, when relevant and proper prior
17 knowledge about the plausibility of alternative conceptual models is included in an analysis
18 following the GLUE-BMA approach, the predictive capacity of the approach is substantially
19 improved.

20

21 **6. Conclusions**

22 We investigated the influence of prior knowledge and prior model probability definition in a
23 multi-model Bayesian averaging methodology which follows Bayesian formalism and that is
24 used to assess uncertainty in the predictions of groundwater models arising from errors in the
25 model structure, input (forcing) data and parameter estimates. The sensitivity analysis was
26 based on the partitioning of the prior model probability space into discrete equidistant
27 intervals of fixed probability. Subsequently, potential combinatorial sets were permuted to

1 obtain sets of prior model probabilities for 7 alternative conceptualizations. The discrete sets
2 were used to numerically analyze the sensitivity of posterior model probabilities and the
3 leading moments of multi-model predictions of groundwater budget terms.

4

5 Additionally, the value of prior knowledge about alternative conceptual models in reducing
6 conceptual model uncertainty was assessed using three illustrative sets of prior model
7 probabilities. The three sets represented knowledge states expressing a non-informative case,
8 proper prior knowledge, and improper prior knowledge about the plausibility of alternative
9 conceptual models. For each of the sets a nonlinear optimization problem was solved in the
10 form of linear (in)equalities expressing quantitative relationships among the alternative
11 conceptualizations. This resulted in three optimized sets of prior model probabilities in
12 agreement with the prior knowledge at hand.

13

14 For illustrative purposes a 3-dimensional hypothetical setup consisting of 2 aquifers separated
15 by an aquitard, in which the flow field was considerably affected by pumping wells and
16 spatially variable hydraulic conductivity, was used. Seven alternative conceptualizations with
17 increasing complexity were adopted to describe the 3-dimensional hypothetical setup. Two of
18 the simpler one-layer models were discarded from further analysis based on the evidence
19 provided by the data.

20

21 Posterior model probabilities and leading moments of the multi-model predictive
22 distributions showed to be very sensitive to different sets of prior model probabilities. This
23 sensitivity clearly states the relevance of selecting proper prior probabilities in the context of
24 the multi-model approach proposed by Rojas *et al.*, (2008). In addition, increasing the prior
25 model probability of a given alternative conceptual model over the other conceptualizations
26 yielded biased leading moments and under-dispersive uncertainty estimations.

1 We showed that an optimized set of prior model probabilities in agreement with proper prior
2 knowledge outperformed the non-informative and improper prior knowledge cases.
3 Reductions between 40 and 60% (with respect to the non-informative case) for the total
4 variances in model predictions were observed when proper prior knowledge was included in
5 the analysis. On the contrary, total variances increased between 32 and 60% respect to the
6 non-informative case when improper prior knowledge was included. Between-model
7 variances, on the other hand, decreased between 50 and 62% when proper prior knowledge
8 was included. Although in absolute terms, between-model and total variances considerably
9 decreased with respect to the non-informative case when proper prior knowledge was
10 included, for the problem at hand, the ratio between-model variance to total variance, within
11 each optimized set, was not substantially modified. This suggests that the contribution of
12 conceptual model uncertainty to total uncertainty can not be further reduced based only on
13 prior knowledge about the plausibility of alternative conceptual models. This implies that
14 other sources of information or conditioning data should be included to further reduce this
15 component of the total variance.

16

17 The results of this study advocate incorporating proper prior knowledge about alternative
18 conceptual models whenever available. Using a 3-dimensional hypothetical setup and three
19 optimized discrete sets of prior model probabilities, it was shown that the predictive
20 performance of the multi-model methodology proposed by Rojas *et al.*, (2008) could be
21 largely improved when proper knowledge is included. It is expected that combining proper
22 prior knowledge about alternative conceptual models with other qualitative or quantitative
23 sources of conditioning data, such as conductivity data, transient groundwater head
24 information or recharge estimates, will further reduce conceptual model uncertainty. These
25 topics will be subject of future research.

1 **Acknowledgments**

- 2 The authors thank the financial support provided to the first author in the framework of the
- 3 PhD IRO-scholarships of the Katholieke Universiteit Leuven (K.U. Leuven). Assistance
- 4 provided by Jorge Gonzalez to implement the R scripts is also acknowledged.

1 **References**

2 Ajami N, Duan Q, Gao X, Sorooshian S. 2005. Multi-model combination techniques for
3 hydrologic forecasting: application to distributed model intercomparison project results.
4 *Journal of Hydrometeorology* **7**(4): 755-768.

5
6 Applebaum D. 1996. *Probability and information: an integrated approach*. Cambridge
7 University Press: New York; 297.

8
9 Beven K, Binley A. 1992. The future of distributed models – model calibration and
10 uncertainty prediction. *Hydrological Processes* **6**(3): 279-283.

11
12 Beven K, Freer J. 2001. Equifinality, data assimilation, and uncertainty estimation in
13 mechanistic modelling of complex environmental systems using the GLUE methodology.
14 *Journal of Hydrology* **249**(1-4): 11-29.

15
16 Beven K. 2005. A manifesto for the equifinality thesis. *Journal of Hydrology* **320**(1-2): 18-
17 36.

18
19 Binley A, Beven K. 2003. Vadose zone flow model uncertainty as conditioned on
20 geophysical data. *Ground Water* **41**(2): 119-127.

21
22 Bredehoeft J. 2003. From models to performance assessment: The conceptualization
23 problem. *Ground Water* **41**(5): 571-577.

24
25 Bredehoeft J. 2005. The conceptualization model problem – surprise. *Hydrogeology Journal*
26 **13**(1): 37-46.

27

1 Carrera J, Alcolea A, Medina A, Hidalgo J, Slooten L. 2005. Inverse problem in
2 hydrogeology. *Hydrogeology Journal* **13**(1): 206-222.
3
4 Deutsch C, Journel A. 1998. *GSLIB Geostatistical software library and user's guide*. Oxford
5 University Press: New York; 384.
6
7 Draper D. 1995. Assessment and propagation of model uncertainty. *Journal of the Royal*
8 *Statistical Society Series B (with discussion)* **57**(1): 45-97.
9
10 Gelman A, Carlin J, Stern H, Rubin D. 2004. *Bayesian data analysis*. Chapman & Hall/CRC:
11 New York; 668.
12
13 Ghosh J, Mohan D, Tapas S. 2006. *An introduction to Bayesian analysis: Theory and*
14 *methods*. Springer texts in Statistics: New York; 352.
15
16 Harbaugh A, Banta E, Hill M, McDonald M. 2000. *MODFLOW-2000, U.S. Geological*
17 *Survey modular ground-water model-user guide to modularization concepts and the ground-*
18 *water flow process*. U.S. Geological Survey. Open File Report, 00-92; 121.
19
20 Harrar W, Sonnenberg T, Henriksen H. 2003. Capture zone, travel time, and solute transport
21 using inverse modelling and different geological models. *Hydrogeology Journal* **11**(5): 536-
22 548.
23
24 Hoeting J, Madigan D, Raftery A, Volinsky C. 1999. Bayesian model averaging: A tutorial.
25 *Statistical Science* **14**(4): 382-417.
26

- 1 Hojberg A, Refsgaard J. 2005. Model uncertainty – parameter uncertainty versus conceptual
2 models. *Water Science & Technology* **52**(6): 177-186.
- 3
- 4 Kass R, Wasserman L. 1996. The selection of prior distributions by formal rules. *Journal of*
5 *the American Statistical Association* **91**(435): 1343-1370.
- 6
- 7 Liang F, Truong Y, Wong W. 2001. Automatic Bayesian model averaging for linear
8 regression and application in Bayesian curve fitting. *Statistica Sinica* **11**(4): 1005-1029.
- 9
- 10 Meyer P, Ye M, Rockhold S, Neuman S, Cantrell K. 2007. *Combined estimation of*
11 *hydrogeologic conceptual model, parameter and scenario uncertainty with application to*
12 *uranium transport at the Hanford site 300 area*. Report NUREG/CR-6940 PNNL-16396, US
13 Nuclear Regulatory Commission, Washington, USA.
- 14
- 15 Neuman S. 2003. Maximum likelihood Bayesian averaging of uncertain model predictions.
16 *Stochastic Environmental Research and Risk Assessment* **17**(5): 291-305.
- 17
- 18 Poeter E, Anderson D. 2005. Multimodel ranking and inference in ground water modelling.
19 *Ground Water* **43**(4): 597-605.
- 20
- 21 Raftery A, Zheng Y. 2003. Discussion: Performance of Bayesian model averaging. *Journal of*
22 *the American Statistical Association* **98**(464): 931-938.
- 23
- 24 Refsgaard J, van der Sluijs J, Brown J, van de Keur P. 2006. A framework for dealing with
25 uncertainty due to model structure error. *Advances in Water Resources* **29**(11): 1586-1597.
- 26

1 Rojas R, Feyen L, Dassargues A. 2008. Conceptual model uncertainty in groundwater
2 modeling: combining generalized likelihood uncertainty estimation and Bayesian model
3 averaging, submitted to *Water Resources Research*.
4
5 Romanowicz R, Beven K, Tawn J. 1994. Evaluation of prediction uncertainty in non-linear
6 hydrological models using a Bayesian approach. In *Statistics for the Environment II; Water
7 Related Issues*, Barnett V, Turkman K (eds.). Wiley: New York; 297-317.
8
9 Rubin Y. 2003. *Applied stochastic hydrogeology*. Oxford University Press: New York; 391.
10
11 Shannon C. 1948. A mathematical theory of communication. *Bell System Technical Journal*
12 **27**: 379-423, 623-656.
13
14 Spellucci P. 1998. An SQP method for general nonlinear programs using only equality
15 constrained subproblems. *Mathematical Programming* **82**(3): 413-448.
16
17 Tamura R. 2007. *Rdonlp2: An R extension library to use Peter Spellucci's DONLP2 from R*.
18 R package version 0.3-1. <http://arumat.ner/Rdonlp2/>.
19
20 Wasserman L. 2000. Bayesian model selection and model averaging. *Journal of*
21 *Mathematical Psychology* **44**: 92-107.
22
23 Ye M, Neuman S, Meyer P. 2004. Maximum likelihood Bayesian averaging of spatial
24 variability models in unsaturated fractured tuff. *Water Resources Research* **40**, W05113,
25 doi:10.1029/2003WR002557.
26

- 1 Ye M, Neuman S, Meyer P, Pohlmann K. 2005. Sensitivity analysis and assessment of prior
2 model probabilities in MLBMA with application to unsaturated fractured tuff. *Water*
3 *Resources Research* **41**, W12429, doi:10.1029/2005WR004260.
- 4
- 5 Ye M, Pohlmann K, Chapman J, Shafer D. 2006. On evaluation of recharge model
6 uncertainty: A priori and a posteriori. In *Proceedings of the International High –level*
7 *Radioactive Waste Management Conference*, Las Vegas, Nevada; 12.

1 **Figures captions**

2 Figure 1: Three-dimensional hypothetical setup including (\odot) observation wells and (\otimes)
3 pumping wells overlain by the groundwater head distribution in the first layer.

4

5 Figure 2: Posterior model probabilities for alternative conceptual models: a) 1Lhtg-L3, b)
6 1Lhtg-AVG, c) 2Lhtg, d) 2LQ3Dhtg and e) 3Lhtg for various sets of discrete prior model
7 probabilities. Symbols represent optimized values of Prior Set 1 (\square), Prior Set 2 (\diamond) and
8 Prior Set 3 (\circ) described in section 4.3.

9

10 Figure 3: Sensitivity analysis in function of prior model probabilities for alternative
11 conceptual model 3Lhtg for: a) prior entropy, b) likelihood ratio (respect to the Prior Set 1)
12 and c) posterior entropy. Symbols represent optimized values of Prior Set 1 (\square), Prior Set 2
13 (\diamond) and Prior Set 3 (\circ) described in section 4.3.

14

15 Figure 4: Leading moments for the posterior predictive distribution of river gains as function
16 of the prior model probabilities of three alternative conceptual models 1Lhtg-L3 (a-d-g),
17 2Lhtg, (b-e-h) and 3Lhtg (c-f-i). Symbols represent optimized values of Prior Set 1 (\square), Prior
18 Set 2 (\diamond) and Prior Set 3 (\circ) described in section 4.3.

19

20 Figure 5: Contours of total variance (a-b-c) and between-model variance (d-e-f) (expressed as
21 a percentage of total variance) for: a) recharge inflows $\times 10^4$ [$\text{m}^3 \text{d}^{-1}$]²; b) river gains $\times 10^4$
22 [$\text{m}^3 \text{d}^{-1}$]², and c) EVT outflows $\times 10^5$ [$\text{m}^3 \text{d}^{-1}$]² in the space of prior model probabilities of
23 alternative conceptual models 1Lhtg-L3 and 3Lhtg when remaining models approach the
24 non-informative case. Symbols represent optimized values of Prior Set 1 (\square), Prior Set 2 (\diamond)
25 and Prior Set 3 (\circ) described in section 4.3.

26

1 Figure 6: Summary statistics for the GLUE-BMA posterior predictive distributions for
2 groundwater budget terms a) WBC inflows, b) recharge inflows, c) WBC outflows, d) river
3 gains and e) EVT outflows for the optimized discrete sets Prior Set 1 (black), Prior Set 2 (red)
4 and Prior Set 3 (light-grey) described in section 4.3. Open circles represent observed values
5 obtained from the 3-dimensional hypothetical setup. Q_1 and Q_3 represent the first and third
6 quartile, respectively. Maximum values are truncated to enhance visual comparison.

7

8 Figure 7: GLUE-BMA posterior mean (diamonds) estimated using equation 3 and the
9 corresponding error bars expressing ± 1 standard deviation (estimated using equation 4) for
10 the sixteen observation wells depicted in Figure 1 for the optimized discrete sets Prior Set 1
11 (black), Prior Set 2 (red) and Prior Set 3 (light-grey) described in section 4.3. Open circles
12 represent observed values obtained from the 3-dimensional hypothetical setup.

1 **Tables**

2 Table 1: Parameters describing the hydraulic conductivity spatial correlation structure for the
 3 different layers of the 3-dimensional hypothetical setup (based on Rubin (2003), Tables 2.1
 4 and 2.2, p34-36).

Layer	Model Parameters		
	μ_K [m d ⁻¹]	σ_{LnK}	I_{LnK}
1	0.1	2.0	400
2	0.01	0.5	800
3	1	1.5	600

5

6

7 Table 2: Range of prior uniform distributions for unknown parameters.

Parameters			Range	
			Minimum	Maximum
Recharge rate	(RECH)	[m d ⁻¹]	0	5.8e-04
Constant head west boundary condition (WBC)	(CH)	[m]	25	75
Elevation surface (EVT)	(SURF)	[m]	30	50
Extinction depth (EVT)	(EXTD)	[m]	0	25
Evapotranspiration rate (EVT)	(EVTR)	[m d ⁻¹]	0	7.0e-03
River conductance	(RIVC)	[m ² d ⁻¹]	1.0e-02	1000

1 Table 3: Summary of integrated model likelihoods and posterior model probabilities for 7 alternative conceptual models and three optimized sets of prior
 2 model probabilities.

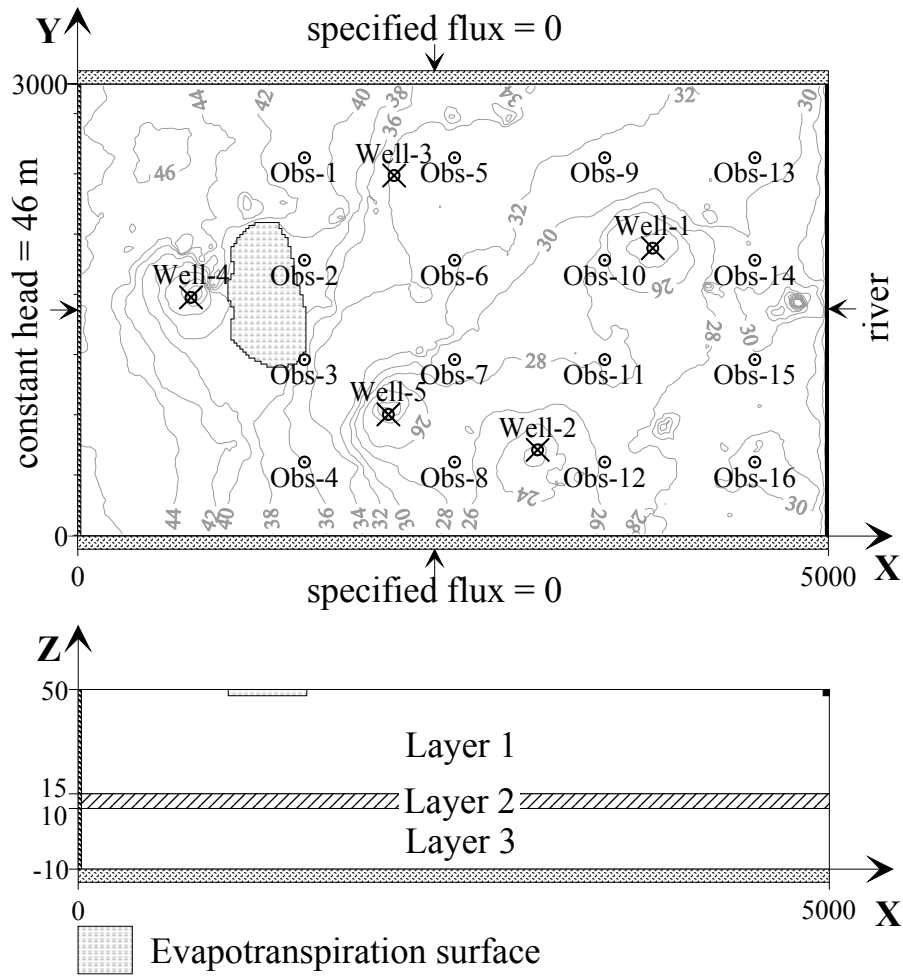
Prior Sets	Conceptual models							Entropy	Likelihood ratio
	1Lhtg-L1	1Lhtg-L2	1Lhtg-L3	1Lhtg-AVG	2Lhtg	2LQ3Dhtg	3Lhtg		
$p(\mathbf{D} \mathbf{M}_k)$	0	0	902.6	935.6	990.4	1046.9	1079.4		
Prior Set 1 $p(\mathbf{M}_k)$	14.29	14.29	14.29	14.29	14.29	14.29	14.29	1.95	1.00
$p(\mathbf{M}_k \mathbf{D})$	0	0	18.22	18.88	19.99	21.13	21.78	1.61	
Prior Set 2 $p(\mathbf{M}_k)$	4.08	4.08	4.08	8.15	12.23	31.79	35.60	1.58	1.34
$p(\mathbf{M}_k \mathbf{D})$	0	0	3.87	8.02	12.73	34.99	40.40	1.32	
Prior Set 3 $p(\mathbf{M}_k)$	23.81	23.81	23.81	11.90	9.52	4.76	2.38	1.74	0.70
$p(\mathbf{M}_k \mathbf{D})$	0	0	43.31	22.45	19.01	10.05	5.18	1.40	

- 1 Table 4: Total variance and between-model variance for groundwater budget terms
- 2 (expressed in $[m^3 d^{-1}]$) as function of the optimized prior probability sets described in section
- 3 4.3. Values in parentheses express percentage reduction with respect to the Prior Set 1.

Groundwater budget terms	Prior Set 1		Prior Set 2		Prior Set 3	
	Total Variance	Between-model variance	Total Variance	Between-model variance	Total Variance	Between-model variance
WBC inflow	463319.2	46854.7	280622.5 (39.4)	17666.0 (62.3)	667218.6 (-44.0)	65680.3 (-40.2)
Recharge inflow	516007.8	86870.7	312683.0 (39.4)	43341.4 (50.1)	681708.6 (-32.1)	79460.8 (8.5)
WBC outflow	7624.7	342.8	3016.7 (60.4)	173.4 (49.4)	12172.6 (-59.6)	319.1 (6.9)
River gains	33893.9	8951.4	19218.0 (43.3)	3871.3 (56.8)	48020.7 (-41.7)	10321.4 (-15.3)
EVT outflow	158321.9	23414.6	82788.2 (47.7)	11495.5 (50.9)	235107.8 (-48.5)	22315.2 (4.7)

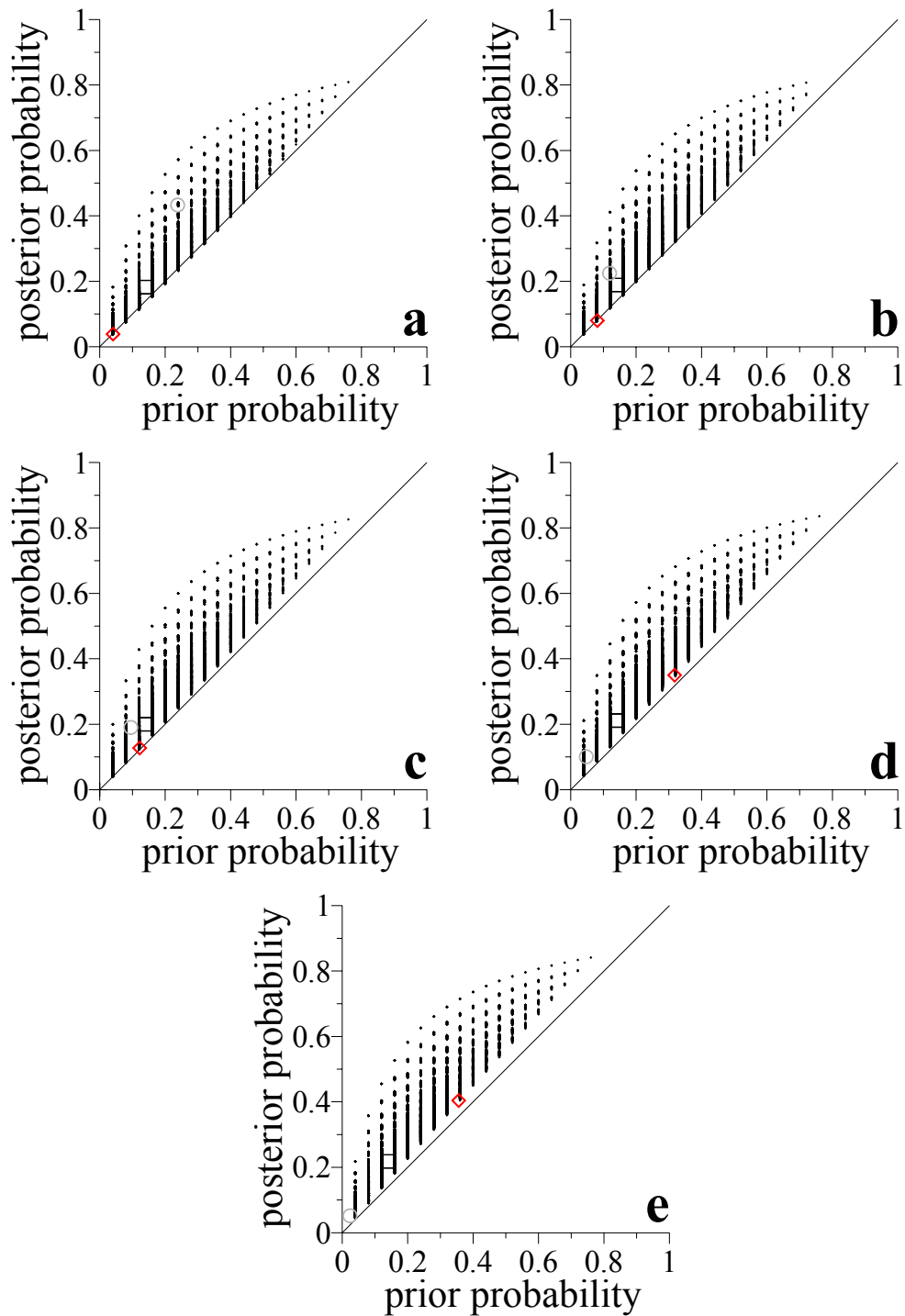
1 **Figures**

2



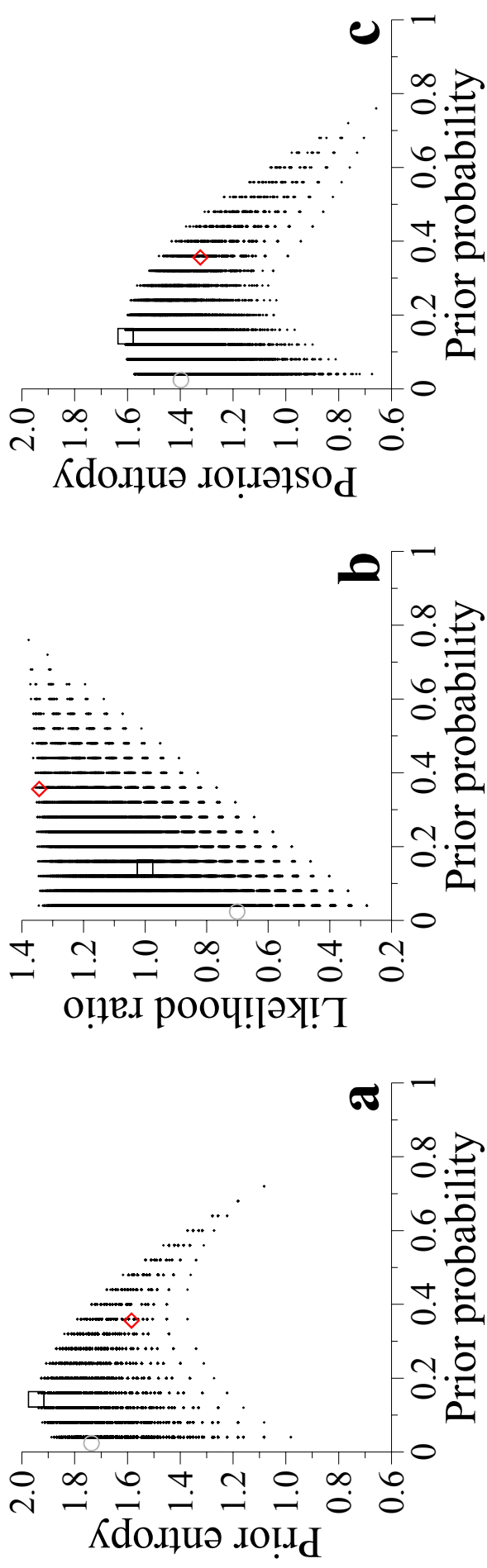
4 Figure 1: Three-dimensional hypothetical setup including (⊙) observation wells and (⊗)

5 pumping wells overlain by the groundwater head distribution in the first layer.

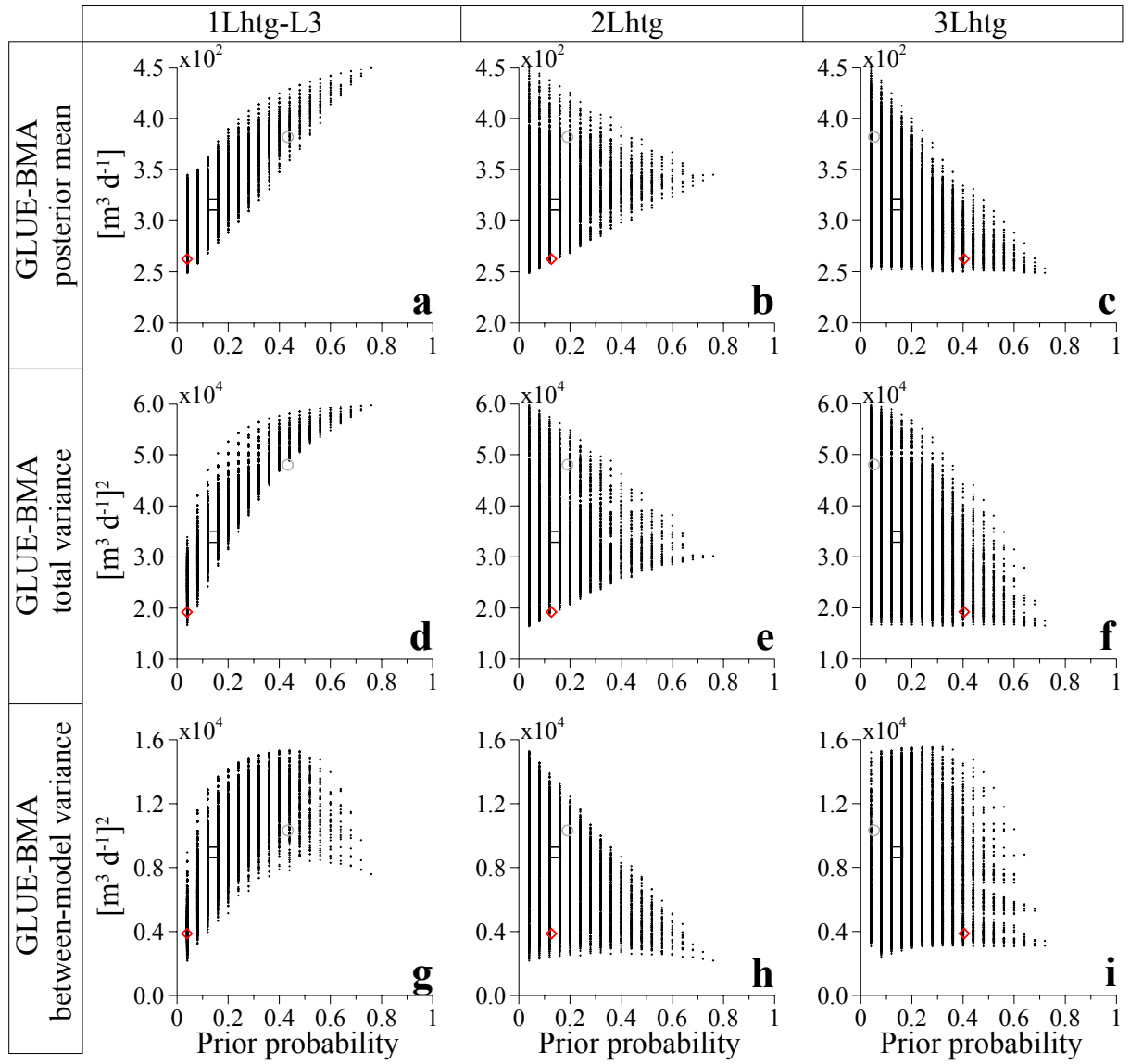


1

2 Figure 2: Posterior model probabilities for alternative conceptual models: a) 1Lhtg-L3, b)
 3 1Lhtg-AVG, c) 2Lhtg, d) 2LQ3Dhtg and e) 3Lhtg for various sets of discrete prior model
 4 probabilities. Symbols represent optimized values of Prior Set 1 (□), Prior Set 2 (◇) and
 5 Prior Set 3 (○) described in section 4.3.

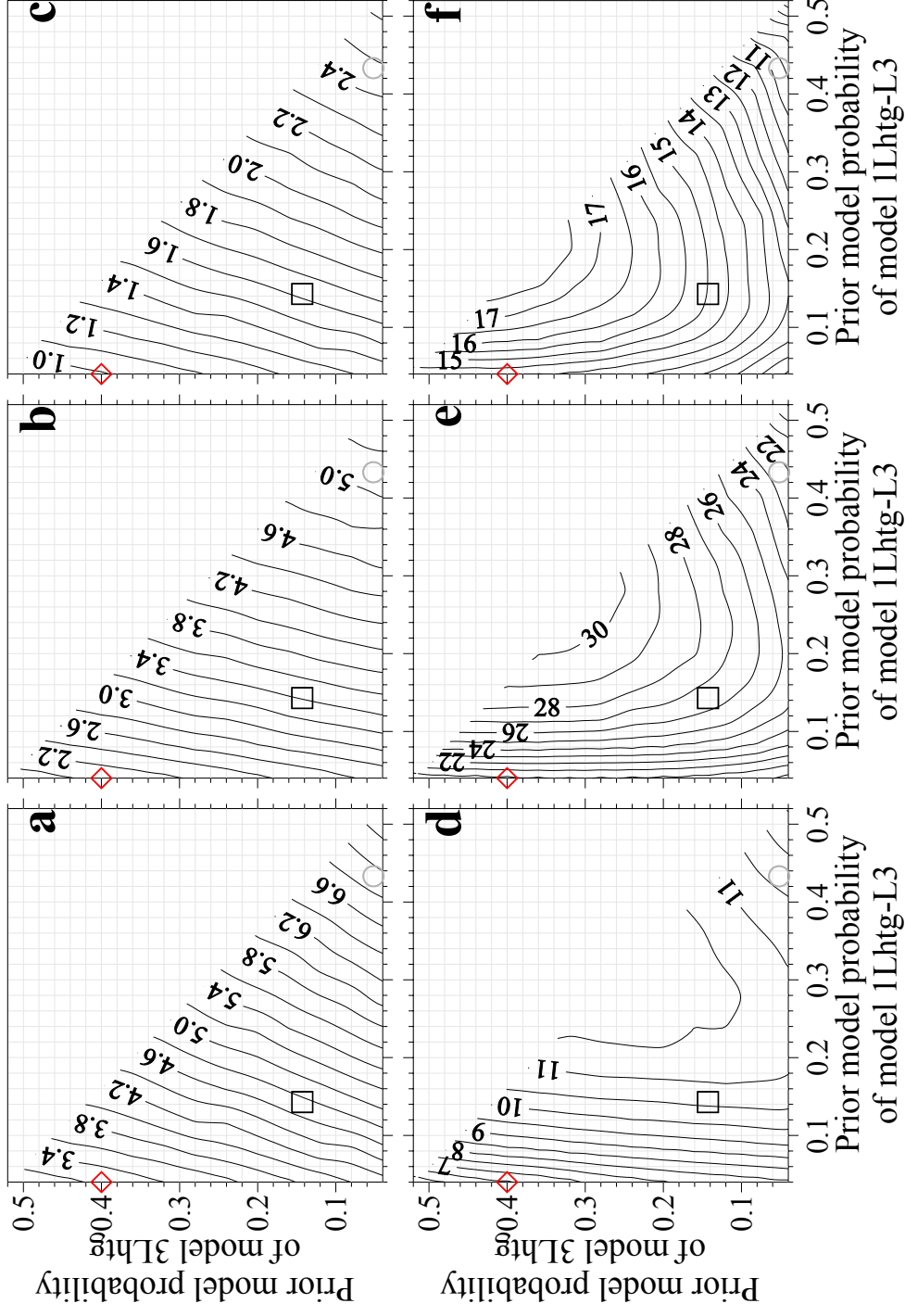


1
2 Figure 3: Sensitivity analysis in function of prior model probabilities for alternative conceptual model 3Lhtg for: a) prior entropy, b) likelihood ratio
3 (respect to the Prior Set 1) and c) posterior entropy. Symbols represent optimized values of Prior Set 1 (\square), Prior Set 2 (\diamond) and Prior Set 3 (\circ) described
4 in section 4.3.

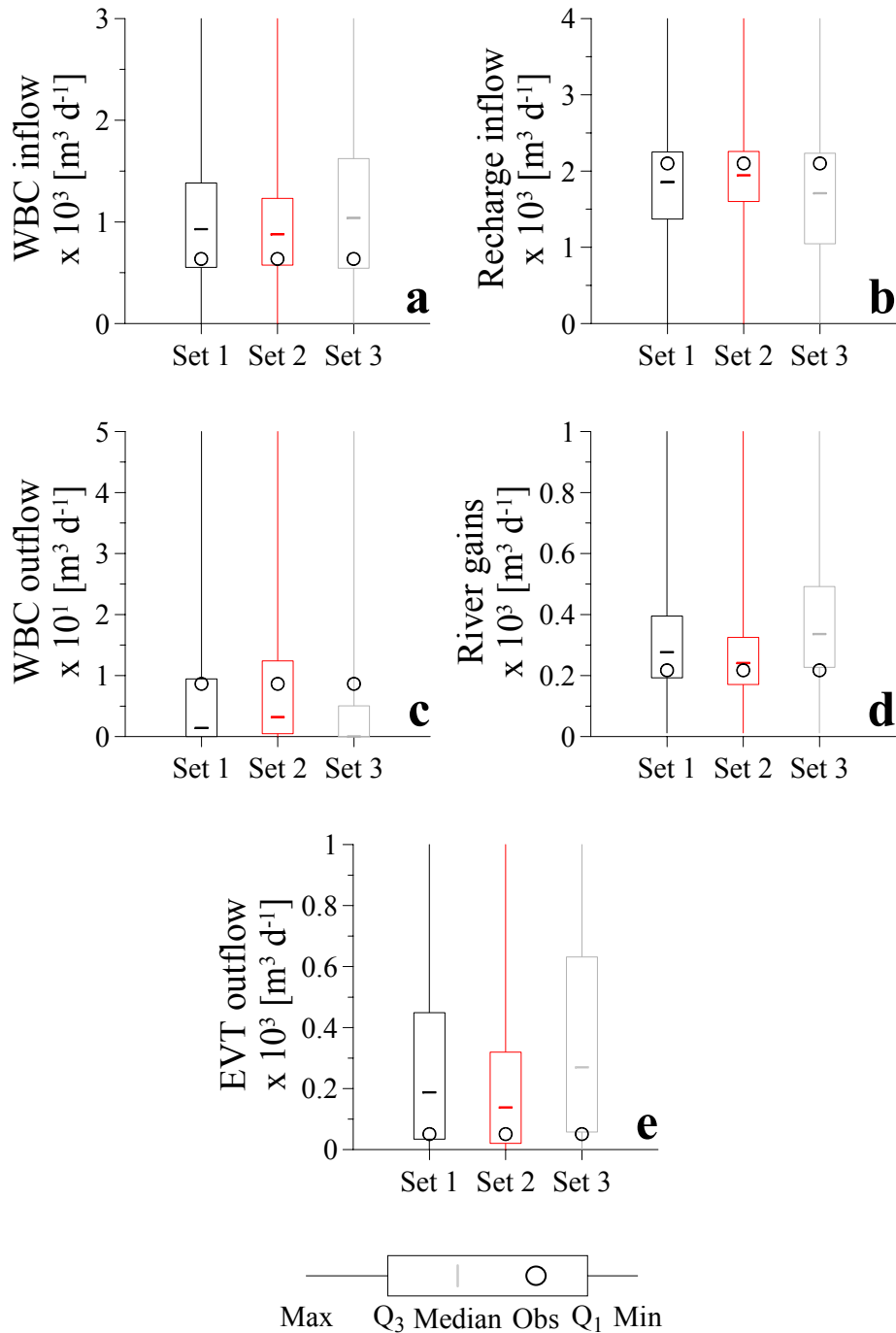


1

2 Figure 4: Leading moments for the posterior predictive distribution of river gains as function
3 of the prior model probabilities of three alternative conceptual models 1Lhtg-L3 (a-d-g),
4 2Lhtg, (b-e-h) and 3Lhtg (c-f-i). Symbols represent optimized values of Prior Set 1 (□), Prior
5 Set 2 (◇) and Prior Set 3 (○) described in section 4.3.

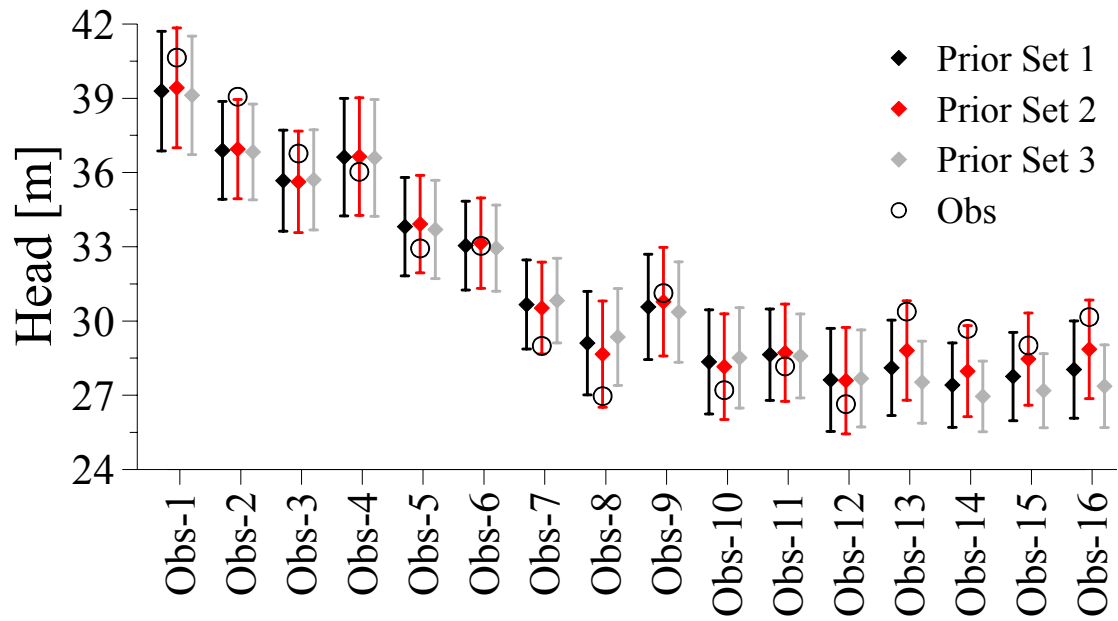


1
2 Figure 5: Contours of total variance (a-b-c) and between-model variance (d-e-f) (expressed as a percentage of total variance) for: a) WBC inflows $\times 10^5$
3 $[\text{m}^3 \text{d}^{-1}]^2$; b) river gains $\times 10^4$ $[\text{m}^3 \text{d}^{-1}]^2$, and c) EVT outflows $\times 10^5$ $[\text{m}^3 \text{d}^{-1}]^2$ in the space of prior model probabilities of alternative conceptual models
4 1Lhtg-L3 and 3Lhtg when remaining models approach the non-informative case. Symbols represent optimized values of Prior Set 1 (\square), Prior Set 2 (\diamond)
5 and Prior Set 3 (\circ) described in section 4.3.



1

2 Figure 6: Summary statistics for the GLUE-BMA posterior predictive distributions for
 3 groundwater budget terms a) WBC inflows, b) recharge inflows, c) WBC outflows, d) river
 4 gains and e) EVT outflows for the optimized discrete sets Prior Set 1 (black), Prior Set 2 (red)
 5 and Prior Set 3 (light-grey) described in section 4.3. Open circles represent observed values
 6 obtained from the 3-dimensional hypothetical setup. Q₁ and Q₃ represent the first and third
 7 quartile, respectively. Maximum values are truncated to enhance visual comparison.



1

2

Figure 7: GLUE-BMA posterior mean (diamonds) estimated using equation 3 and the corresponding error bars expressing ± 1 standard deviation (estimated using equation 4) for the sixteen observation wells depicted in Figure 1 for the optimized discrete sets Prior Set 1

3

(black), Prior Set 2 (red) and Prior Set 3 (light-grey) described in section 4.3. Open circles

4

represent observed values obtained from the 3-dimensional hypothetical setup.

5

6