

1 **Conceptual model uncertainty in groundwater modeling: combining**
2 **generalized likelihood uncertainty estimation and Bayesian model**
3 **averaging**

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

6
7 ^{a*} Corresponding author.

8 Applied geology and mineralogy, Department of Earth and Environmental Sciences,
9 Katholieke Universiteit Leuven, Celestijnenlaan 200 E, 3001 Heverlee, Belgium

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

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

12
13 ^b Land management and natural hazards unit, Institute for Environment and Sustainability,
14 DG Joint Research Centre, European Commission, TP261, 21020 Ispra (Va), Italy

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

16
17 ^c Hydrogeology and environmental geology, Department of Architecture, Geology,
18 Environment, and Civil Engineering (ArGEnCo),

19 Université de Liège, B.52/3 Sart-Tilman, B-4000 Liège, Belgium

20 E-mail address: alain.dassargues@geo.kuleuven.be, alain.dassargues@ulg.ac.be

1 **Abstract**

2 Uncertainty assessments in groundwater modeling applications typically attribute all sources
3 of uncertainty to errors in parameters and inputs, neglecting what may be the primary source
4 of uncertainty, namely, errors in the conceptualization of the system. Confining the set of
5 plausible system representations to a single model leads to under-dispersive and prone to bias
6 predictions. In this work we present a general and flexible approach that combines
7 Generalized Likelihood Uncertainty Estimation (GLUE) and Bayesian Model Averaging
8 (BMA) to assess uncertainty in model predictions that arise from errors in model structure,
9 inputs and parameters. In a prior analysis, a set of plausible models are selected and the joint
10 prior input and parameter space is sampled to form potential simulators of the system. For
11 each model the likelihood measures of acceptable simulators, assigned to them based on their
12 ability to reproduce observed system behavior, are integrated over the joint input and
13 parameter space to obtain the integrated model likelihood. The latter is used to weight the
14 predictions of the respective model in the BMA ensemble predictions. For illustrative
15 purposes we applied the methodology to a three-dimensional hypothetical setup. Results
16 showed that predictions of groundwater budget terms varied considerably among competing
17 models, although that a set of 16 head observations used for conditioning did not allow
18 differentiating between the models. BMA provided consensus predictions that were more
19 conservative. Conceptual model uncertainty contributed up to 30% of the total uncertainty.
20 The results clearly indicate the need to consider alternative conceptualizations to account for
21 model uncertainty.

22

23 Keywords: GLUE, BMA, Multi-model prediction, Monte Carlo methods, uncertainty
24 assessment

1 **1. Introduction and scope**

2 With increasing human and climate pressures on groundwater resources, accurate and reliable
3 predictions of groundwater flow and pollutant transport are essential for sustainable
4 groundwater management practices. However, typically, the geological structure is only
5 partially known and point measurements of subsurface properties or groundwater heads are
6 sparse and prone to error. Consequently, incomplete or biased process representation, errors
7 in the specification of initial and boundary conditions, as well as errors in the model
8 parameters, render the predictions of groundwater dynamics and pollutant transport uncertain.

9
10 Over the last decades, considerable efforts have been put in developing methods to determine
11 optimal groundwater parameter values and in quantifying model prediction uncertainty
12 associated with uncertainty in these parameter estimates. This has resulted in a variety of
13 inverse techniques for groundwater modeling applications. We do not wish to provide a
14 complete overview of parameter estimation methods but refer the reader to Sun (1994) and
15 Carrera et al. (2005) for excellent reviews. Despite its extensive application, the major
16 weakness of parameter-calibration approaches is that all sources of uncertainty are attributed
17 to parameter errors. This often results in biased parameter estimates that compensate for
18 errors in model structure, input data and measurement errors.

19
20 Typically, these methods ignore conceptual or structural uncertainty by confining the range
21 of plausible system representations to a single hydrological model. This often leads to
22 overconfidence in the predictive capabilities of the model and in predictive uncertainty
23 analyses that are under-dispersive and prone to statistical bias. In recent years, a number of
24 authors have acknowledged that conceptual model uncertainty has received less formal
25 attention in groundwater applications than it should (e.g., Neuman, 2003, Neuman and
26 Wierenga, 2003; Bredehoeft, 2003, 2005; Carrera et al., 2005; Poeter and Anderson, 2005;
27 Refsgaard et al., 2006). Bredehoeft (2005) summarizes the main issues concerning conceptual

1 model specification as follows: (1) modelers tend to consider their conceptual models as
2 immutable; (2) frequently, errors in model predictions turn around a poor choice of the
3 conceptual model; (3) data will fit more than one conceptual model equally well;
4 consequently, (4) a good calibration of a model does not ensure a correct conceptual model;
5 and (5) parametric uncertainty does not compensate for uncertainties derived from the
6 conceptual model specification.

7

8 These concerns have motivated researchers in the hydrological sciences to consider multi-
9 model methods, which seek to obtain consensus predictions from a set of plausible models by
10 linearly combining individual model predictions. The weights to aggregate multiple model
11 outputs can be equal (model average) in the simplest case, or can be determined through
12 regression-based approaches (e.g., Abrahart and See, 2002; Georgakakos et al., 2004).
13 However, the weights in such combinations are not connected to model performance and can
14 take any arbitrary value, hence lacking physical interpretation. An approach in which weights
15 are intrinsically connected to model performance has been proposed by Poeter and Anderson
16 (2005). This approach combines predictions of multiple competing models using Akaike's
17 weights (Akaike, 1974; Burnham and Anderson, 1998). However, it lacks of a consistent way
18 to incorporate previous knowledge about parameters and conceptual models in the multi-
19 model prediction. A similar method that partially overcomes the restriction of including
20 previous knowledge about multiple model structures has been proposed by Refsgaard et al.,
21 (2006). In this approach, a suite of conceptual models are independently calibrated and a
22 pedigree analysis is performed to assess the overall tenability of the models. Nonetheless, the
23 pedigree analysis does not provide an indication of the relative quality of the different model
24 structures and, consequently, it is difficult to include it in a quantitative uncertainty analysis
25 in terms of probabilities (Refsgaard et al., 2006).

1 Bayesian Model Averaging (BMA) (Draper, 1995; Hoeting et al., 1999), on the other hand,
2 employs probabilistic techniques to derive consensus predictions from a set of alternative
3 models. In short, BMA weights the predictions of competing models by their corresponding
4 posterior model probability, representing each model's relative skill to reproduce system
5 behavior in the training period. Hence, BMA weights are tied directly to individual model
6 performance. Several studies applying the method to a range of different problems have
7 demonstrated that BMA produces more accurate and reliable predictions than other existing
8 multi-model techniques (e.g., Raftery and Zheng, 2003; Ye et al., 2004; Ajami et al., 2005).
9
10 In the field of groundwater hydrology applications of BMA have been rare. Neuman (2003)
11 proposed the Maximum Likelihood Bayesian Model Averaging (MLBMA) method to assess
12 the joint predictive distribution of several competing models. MLBMA is an approximation
13 of BMA that relies on maximum likelihood parameter estimation and expanding around these
14 values through Monte Carlo simulation. Subsequently, the posterior model probabilities are
15 approximated using the Kashyap information criterion (Kashyap, 1982). MLBMA does not
16 require exhaustive Monte Carlo simulations and obviates the need of (though it can
17 incorporate) prior information about model parameters, which is often difficult to obtain (Ye
18 et al., 2005). Ye et al. (2004) expanded upon the theoretical framework of MLBMA and
19 applied it to model the log permeability in unsaturated fractured tuff using alternative
20 variogram models.
21
22 An alternative methodology that rejects the idea of a unique optimal simulator of the natural
23 system is the Generalized Likelihood Uncertainty Estimation (GLUE) method (Beven and
24 Binley, 1992; Beven, 1993). GLUE is based on the concept of equifinality, which
25 acknowledges that there exist many combinations of model structures and parameter sets that
26 provide (equally) good reproductions of the observed system response. For each possible
27 simulator a likelihood measure is defined based on the degree of correspondence between

1 simulated and observed records of system responses. Simulators that perform better than a
2 subjectively chosen threshold criteria are retained and consequently used to provide an
3 ensemble of likelihood weighted predictions of the system under future forcing conditions.
4 The technique has found its application mainly in rainfall-runoff and flood inundation
5 modeling (see e.g., Beven and Freer (2001) and Beven (2005b) for a complete list of
6 references). In recent years, GLUE has also been applied in several groundwater studies (e.g.,
7 Feyen et al., 2001; Binley and Beven, 2003; Morse et al., 2003). Even though equifinality, as
8 defined by Beven (1993; 2005a), arises because of the combined effects of errors in the
9 forcing data, system conceptualization, measurements and parameter estimates, as yet, it has
10 only been applied in the context of a single deterministic conceptual model (Refsgaard et al.,
11 2006), thereby, neglecting model structural uncertainty.

12

13 In this work, we combine GLUE with BMA to explicitly account for uncertainty that
14 originates from errors in the model conceptualization, forcing data (e.g., recharge rate,
15 boundary conditions) and parameter values. Within the GLUE framework, we explore the
16 global likelihood response surface of all possible combinations of plausible model structures,
17 forcing data and parameter values in order to select those simulators that perform well. For
18 each model structure, the posterior model probability is obtained by integrating the likelihood
19 measures over the retained simulators for that model structure. The posterior model
20 probabilities are subsequently used in BMA to weight the predictions of the competing
21 models when assessing the joint predictive uncertainty.

22

23 The method presented is very flexible since (i) there is no restriction on the diversity of
24 conceptual models or on the level of uncertainty in the forcing data or parameters that can be
25 included; (ii) it allows for different ways of expressing the likelihood of a simulator
26 (including a formal Bayesian one) based on the distribution of the residuals, hence allowing
27 different types of knowledge to be incorporated (quantitative as well as qualitative); and (iii)

1 it is Bayesian in nature, which provides a formal framework to incorporate previous
2 knowledge about the model structures and parameters, or to update the estimates should new
3 information become available. The main drawback of the methodology is the computational
4 burden. Due to the presence of multiple local optima in the global likelihood response
5 surface, good performing or behavioral simulators might be well distributed across the
6 hyperspace dimensioned by the set of model structures, input and parameter vectors. This
7 necessitates that the global likelihood surface is extensively sampled.

8

9 The remainder of this paper is organized as follows. In section 2, we provide a condensed
10 overview of the GLUE and BMA methodologies, followed by a description of the procedure
11 to integrate both methods. Section 3 details a three-dimensional hypothetical setup that is
12 used to illustrate the integrated uncertainty assessment methodology. Implementation details
13 are described in section 4. In this section we elaborate on the different conceptualizations as
14 well as on input and parameter uncertainty. Results are discussed in section 5 and a summary
15 of conclusions is presented in section 6.

16

17 **2. Methodology for integrated uncertainty assessment**

18 **2.1. Generalized Likelihood Uncertainty Estimation (GLUE) methodology**

19 GLUE is a Bayesian Monte Carlo simulation technique based on the concept of equifinality
20 (Beven and Binley, 1992; Beven and Freer, 2001). It rejects the idea of a single correct
21 representation of the system in favor of many acceptable or behavioral system representations
22 that should be considered in the evaluation of uncertainty associated with predictions (Beven,
23 2005b). For each simulator sampled from a prior set of possible system representations a
24 likelihood measure is calculated that reflects the ability of the simulator to simulate the
25 system responses, given the available training data. Simulators that perform below a rejection
26 criterion are discarded from the further analysis and the likelihood measures of retained
27 simulators are rescaled so as to render the cumulative likelihood equal to one. Ensemble

1 predictions are based on the predictions of the retained set of simulators, weighted by their
2 respective rescaled likelihood.

3

4 The likelihood or “goodness of fit” used in GLUE must be seen in a much wider sense than
5 the formal likelihood functions used in traditional statistical estimation theory. The
6 likelihoods used in GLUE are a measure of the ability (performance) of a simulator to
7 reproduce a given set of training data. Therefore, they represent an expression of belief in the
8 predictions of that particular simulator rather than a formal definition of probability being the
9 correct representation of the system (Binley and Beven, 2003). However, the GLUE
10 methodology is fully coherent with a formal Bayesian approach when the use of a classical
11 likelihood function is justifiable based on the nature of the residuals (see e.g., Romanowicz et
12 al., 1994).

13

14 Some critiques have recently been raised concerning the subjective nature of some decisions
15 that have to be made in order to implement the GLUE methodology (see e.g., Mantovan and
16 Todini (2006) and the reply of Beven et al., (2007)). These subjective decisions involve the
17 definition of a suitable likelihood function and the definition of the rejection level in order to
18 distinguish between “behavioral” and “non-behavioral” simulators. To evaluate the impact of
19 these subjective decisions in the analysis, we implement three different likelihood functions
20 in this study, namely, a formal statistical, a GLUE type, and a Fuzzy type measure.

21

22 Let us consider a set of plausible model structures $\mathbf{M}=\{M_1, M_2, \dots, M_k, \dots, M_K \mid K < \infty\}$, a set
23 of parameter vectors $\Theta=(\theta_1, \theta_2, \dots, \theta_l, \dots, \theta_L)$ and a set of input variable vectors
24 $\Upsilon=(Y_1, Y_2, \dots, Y_m, \dots, Y_M)$, and denote the observed and simulated system variable vectors as
25 $\mathbf{D}=(D_1, D_2, \dots, D_n, \dots, D_N)$ and $\mathbf{D}^*=(D_1^*, D_2^*, \dots, D_n^*, \dots, D_N^*)$, respectively. Then,

1 $L(\mathbf{M}_k, \boldsymbol{\theta}_l, \mathbf{Y}_m | \mathbf{D})$ represents the likelihood of the k th model structure parameterized with
 2 parameter vector $\boldsymbol{\theta}_l$ and forced by input data vector \mathbf{Y}_m to represent the true system, given
 3 the observations in \mathbf{D} .

4

5 As a first likelihood measure, we consider a Gaussian likelihood function (1), which is based
 6 on the assumption that the residuals follow a normal distribution centered on zero. For a
 7 given number of observations, N , the Gaussian likelihood is given by

8

$$9 \quad L(\mathbf{M}_k, \boldsymbol{\theta}_l, \mathbf{Y}_m | \mathbf{D}) = (2\pi)^{-N/2} |C_{\mathbf{D}}|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{D} - \mathbf{D}^*)^T C_{\mathbf{D}}^{-1}(\mathbf{D} - \mathbf{D}^*)\right) \quad (1)$$

10

11 where, $C_{\mathbf{D}}$ is the covariance matrix of the observed system variables.

12

13 The second measure implemented is the model efficiency likelihood function (2) (Freer and
 14 Beven, 1996; Feyen et al., 2001; Jensen, 2003), which is based on the Nash-Sutcliffe
 15 efficiency criterion (Nash and Sutcliffe, 1970) with shaping factor S , and is given by

16

$$17 \quad L(\mathbf{M}_k, \boldsymbol{\theta}_l, \mathbf{Y}_m | \mathbf{D}) = \left(1 - \frac{\sigma_{\varepsilon}^2}{\sigma_{\mathbf{D}}^2}\right)^S \quad (2)$$

18

19 where σ_{ε}^2 and $\sigma_{\mathbf{D}}^2$ are the variance of the residuals and of the observations, respectively. We
 20 used a shaping factor S equal to one, in which case the model efficiency likelihood function is
 21 equivalent to the coefficient of determination (R^2).

22

23 As a third measure we implemented a triangular likelihood function (3), belonging to the so-
 24 called Fuzzy type measures (Jensen, 2003), given by

1

$$2 \quad L_n(\mathbf{M}_k, \boldsymbol{\theta}_l, \mathbf{Y}_m | \mathbf{D}_n) = \frac{D_n^* - a}{b - a} I_{a,b} + \frac{c - D_n^*}{c - b} I_{b,c} \quad (3)$$

3

4 where

$$5 \quad I_{a,b} = \begin{cases} 1 & \text{if } a < D_n^* \leq b \\ 0 & \text{otherwise} \end{cases}$$
$$6 \quad I_{b,c} = \begin{cases} 1 & \text{if } b < D_n^* < c \\ 0 & \text{otherwise} \end{cases}$$

6

7 and the limits a and c define the tolerable error. The triangular likelihood function in (3)

8 gives a point likelihood measure $L_n(\mathbf{M}_k, \boldsymbol{\theta}_l, \mathbf{Y}_m | \mathbf{D}_n)$ for each observation data n . It was

9 combined by a geometric mean inference function to obtain a global likelihood value

$$10 \quad L(\mathbf{M}_k, \boldsymbol{\theta}_l, \mathbf{Y}_m | \mathbf{D}).$$

11

12 **2.2. Bayesian Model Averaging (BMA)**

13 BMA provides a coherent framework for combining predictions from multiple competing

14 conceptual models to provide a more realistic and reliable description of the total prediction

15 uncertainty. It is a statistical procedure that infers consensus predictions by weighing

16 predictions from competing models based on their relative skill, with predictions from better

17 performing models receiving higher weights than those of worse performing models.

18

19 Following the notation of Hoeting et al., (1999), if Δ is a quantity to be predicted, the BMA

20 predictive distribution of Δ is given by

21

$$22 \quad p(\Delta | \mathbf{D}) = \sum_{k=1}^K p(\Delta | \mathbf{D}, \mathbf{M}_k) p(\mathbf{M}_k | \mathbf{D}) \quad (4)$$

1

2 Equation (4) is an average of the predictive distributions of Δ under each model considered,3 $p(\Delta | \mathbf{D}, M_k)$, weighted by its posterior model probability, $p(M_k | \mathbf{D})$. This latter term4 reflects how well model k fits the observed data \mathbf{D} and can be computed using Bayes' rule

5

$$6 \quad p(M_k | \mathbf{D}) = \frac{p(\mathbf{D} | M_k) p(M_k)}{\sum_{k'=1}^K p(\mathbf{D} | M_{k'}) p(M_{k'})} \quad (5)$$

7

8 where $p(M_k)$ is the prior probability of model M_k , and $p(\mathbf{D} | M_k)$ is the integrated9 likelihood of model M_k given by

10

$$11 \quad p(\mathbf{D} | M_k) = \iint p(\mathbf{D} | M_k, \boldsymbol{\theta}_l, \mathbf{Y}_m) p(\boldsymbol{\theta}_l, \mathbf{Y}_m | M_k) d\boldsymbol{\theta}_l d\mathbf{Y}_m \quad (6)$$

12

13 where $p(\mathbf{D} | M_k, \boldsymbol{\theta}_l, \mathbf{Y}_m)$ is the likelihood of model structure M_k parameterized with14 parameter vector $\boldsymbol{\theta}_l$ and forced by input data vector \mathbf{Y}_m given the observations in \mathbf{D} , and15 $p(\boldsymbol{\theta}_l, \mathbf{Y}_m | M_k)$ is the joint prior probability distribution of $(\boldsymbol{\theta}_l, \mathbf{Y}_m)$ given model M_k .

16

17 The leading moments of the BMA prediction of Δ are given by (Draper, 1995)

18

$$19 \quad E[\Delta | \mathbf{D}] = \sum_{k=1}^K E[\Delta | \mathbf{D}, M_k] p(M_k | \mathbf{D}) \quad (7)$$

20

$$21 \quad \begin{aligned} Var[\Delta | \mathbf{D}] = & \sum_{k=1}^K Var[\Delta | \mathbf{D}, M_k] p(M_k | \mathbf{D}) \\ & + \sum_{k=1}^K (E[\Delta | \mathbf{D}, M_k] - E[\Delta | \mathbf{D}])^2 p(M_k | \mathbf{D}) \end{aligned} \quad (8)$$

1

2 In essence, the BMA prediction is the weighted average of predictions from a suite of
3 alternative models, with the weights equal to the likelihood that a model represents the true
4 unknown model. From equation (8) it is seen that the variance of the BMA predictions
5 consists of two terms, the first representing the within-model variance and the second
6 representing the between-model variance.

7

8 **2.3. Combining GLUE and BMA**

9 Combining the GLUE and BMA methods involves the following sequence of steps

- 10 1. On the basis of prior and expert knowledge about the site, a suite of alternative
11 conceptualizations is proposed, following, for instance, the methodology proposed by
12 Neuman and Wierenga (2003).
- 13 2. Realistic prior ranges are defined for the input and parameter vectors under each plausible
14 model structure.
- 15 3. A likelihood measure and rejection criteria are defined.
- 16 4. For the suite of alternative conceptual models, input and parameter values are sampled
17 from the prior ranges to generate possible representations or simulators of the system.
- 18 5. A likelihood measure is calculated for each simulator based on the agreement between the
19 simulated and observed system response.
- 20 6. Simulators that are not in agreement with the selected rejection criterion are discarded
21 from the analysis by setting their likelihood to zero.
- 22 7. For each conceptual model M_k a subset A_k of simulators with likelihood
23 $p(\mathbf{D}|M_k, \boldsymbol{\theta}_l, \mathbf{Y}_m) = L(M_k, \boldsymbol{\theta}_l, \mathbf{Y}_m | \mathbf{D})$ is retained. Steps 4-6 are repeated until the
24 hyperspace of possible simulators is adequately sampled, i.e., when the conditional
25 distributions of predicted state variables based on the likelihood weighted simulators in
26 the subset A_k converge to stable distributions for each of the conceptual models M_k .

1 8. The integrated likelihood of each conceptual model M_k (equation 6) is approximated by
 2 summing the likelihood weights of the retained simulators in subset A_k , or

$$3$$

$$4 \quad p(\mathbf{D}|M_k) \approx \sum_{l,m \in A_k} L(M_k, \boldsymbol{\theta}_l, \mathbf{Y}_m | \mathbf{D}) \quad (9)$$

5

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

$$8$$

$$9 \quad p(M_k | \mathbf{D}) \approx \frac{\sum_{A_k} L(M_k, \boldsymbol{\theta}_l, \mathbf{Y}_m | \mathbf{D}) p(M_k)}{\sum_{j=1}^K \sum_{l,m \in A_j} L(M_j, \boldsymbol{\theta}_l, \mathbf{Y}_m | \mathbf{D}) p(M_j)} \quad (10)$$

10

11 10. After normalization of the likelihood weighted predictions under each individual model
 12 (such that the cumulative likelihood under each model equals one) a multi-model
 13 prediction is obtained with equation (4) using the weights obtained in (10). The leading
 14 moments of this distribution are obtained with equations (7) and (8).

15

16 Details about the implementation of the methodology, applied to the three-dimensional
 17 hypothetical setup described in the next section, are presented in Section 4.

18

19 **3. Three-dimensional hypothetical case**

20 For illustrative purposes, we employ a hypothetical setup for which the true conditions are
 21 known. The three-dimensional example system is similar to the reference case described in

22 Poeter and Anderson (2005) and is presented in Figure 1. Lateral dimensions are 5000 m (E-

23 W) by 3000 m (N-S) discretized in 25 m by 25 m grid cells. The system extends over 60 m in

24 the vertical direction, with undisturbed layer thicknesses of 35 m (upper aquifer), 5 m (middle

1 aquitard) and 20 m (lower aquifer). We assume statistically homogeneous deposits with a
2 constant mean hydraulic conductivity K (see Table 1). Smaller-scale variability is represented
3 using the theory of random space functions, adopting isotropic exponential covariance
4 functions for $\ln K$ in all layers. The spatial distribution of the hydraulic conductivity in the
5 layers of the example setup, as well as any other realization of the hydraulic conductivity
6 field used in this work, is generated using the sequential Gaussian simulation (sgsim)
7 algorithm of the Geostatistical Software Library (Deutsch and Journel, 1998). Parameters of
8 the covariance function of $\ln K$ for the different layers are presented in Table 1.

9

10 Simulation of steady-state flow employs Modflow-2000 (Harbaugh et al., 2000). At the north
11 and south boundaries, as well as at the bottom of the lower layer, zero gradient conditions are
12 imposed. A uniform recharge of $1.4 \times 10^{-4} \text{ m d}^{-1}$ is applied to the top layer. At the west
13 boundary a constant head $h = 46 \text{ m}$ is defined. The east side of the domain is bounded by a 10
14 m-wide river with a constant stage of 25 m. The river bottom is at 20 m, defining a constant
15 river water depth of 5 m. It is underlain by 5 m-thick sediments with a vertical hydraulic
16 conductivity of 0.1 m d^{-1} . Five pumping wells are distributed in the area and pump a total of
17 $2450 \text{ m}^3 \text{ d}^{-1}$ from the lowermost layer (Figure 1). An evapotranspiration zone, delineated by
18 the polygon in Figure 1, is defined with an evapotranspiration surface elevation at 43 m, an
19 evapotranspiration rate of $1.37 \times 10^{-3} \text{ m d}^{-1}$ and an extinction depth of 5 m.

20

21 The resulting “true” groundwater head distribution for the top layer is presented as an overlay
22 in Figure 1. The ambient background gradient from west to east is altered considerable by the
23 cones of depression around the pumping wells, local effects of spatially varying hydraulic
24 conductivity and to a lesser extent by the evapotranspiration zone. From the “true”
25 groundwater head distribution for layer 1, values are selected at the 16 locations defined by
26 the observation wells in Figure 1, which are used to estimate the likelihood weights in the
27 evaluation of different simulators.

1 **4. Implementation of the methodology**

2 **4.1. Alternative conceptual models**

3 Theoretically, all possible models of relevance could be included in \mathbf{M} . However, the
4 number of potentially feasible models may be exceedingly large, rendering their exhaustive
5 inclusion in \mathbf{M} infeasible (Hoeting et al., 1999). We adopt the idea of Ockham's Window
6 (Madigan and Raftery, 1994) to consider a relatively small set of the most parsimonious
7 models in view of the knowledge about the system and their ability to explain the data. As a
8 consequence, the joint predictions do not represent all possibilities but only a limited range,
9 conditional on the ensemble of conceptual models used to describe the groundwater system.

10

11 We consider the following seven conceptualizations with increasing complexity to describe
12 the three-dimensional hypothetical setup presented in section 3: (1), (2) and (3) one-layer
13 models with mean K and spatial correlation law of layer 1 (1Lhtg-L1), layer 2 (1Lhtg-L2) and
14 layer 3 (1Lhtg-L3) of the three-dimensional hypothetical setup, respectively; (4) a one-layer
15 model with average mean K and spatial correlation (1Lhtg-AVG); (5) a two-layer model with
16 mean K and spatial correlation taken from layer 1 and layer 3 (2Lhtg); (6) a two-layer quasi-
17 three-dimensional model with mean K and spatial correlation taken from layer 1 and layer 3,
18 and mean K of layer 2 used to define the aquitard (2LQ3Dhtg); and (7) a three-layer model
19 based on the spatial K distributions of layer 1, layer 2 and layer 3 (3Lhtg). All
20 conceptualizations comprise a total aquifer thickness of 60 m and are forced by identical
21 types of boundary conditions, although the magnitudes or rates of the latter are set variable
22 (see 4.3).

23

24 **4.2. Parameterization**

25 The focus of this work is on the assessment of conceptual uncertainty. Therefore, we confine
26 the dimensionality of the analysis by considering uncertainty only in the input variables and
27 parameters related to the evapotranspiration process, lateral boundary conditions, river

1 description and recharge process, i.e., input variables and parameters that are common to all
2 setups (see Table 2). Realizations of the hydraulic conductivity field of the different layers in
3 the alternative conceptualizations are generated with the same mean K and spatial correlation
4 law as the respective layers in the three-dimensional hypothetical setup (values listed in Table
5 1). For the 1Lhtg-AVG conceptualization the averages of these values are used. Although the
6 nature of the underlying structure is assumed to be known, hence only the realization space is
7 sampled, uncertainty in the mean hydraulic conductivity and spatial correlation function can
8 be accounted for using, for example, the Bayesian methods presented in Feyen et al., (2002).

9

10 **4.3. Prior distributions**

11 We assign equal prior probabilities to the seven conceptualizations and adopt uniform prior
12 distributions for the unknown inputs and parameters. The definition of such non-informative
13 prior distributions is based on what is known as the principle of insufficient reason or the
14 Bayes-Laplace postulate. According to this principle, in the absence of evidence to the
15 contrary, all possibilities should have the same initial probability (Bernardo and Smith,
16 2000). Using these priors, we expect that the information in the data, expressed by the
17 likelihood function, should dominate the form of the resulting posterior distribution. The
18 ranges that describe the prior uniform distributions of the unknown variables are presented in
19 Table 2.

20

21 **4.4. Simulation**

22 Parameter and input vectors, sampled from the prior distributions using a Latin Hypercube
23 Sampling (LHS) scheme, are combined with hydraulic conductivity realizations for the
24 respective layers and consequently evaluated under each conceptual model. On the basis of
25 the evaluation of a set of initial runs, a rejection threshold is defined corresponding to a
26 maximum allowable deviation of 5 m at any of the 16 observation wells depicted in Figure 1.
27 A point rejection threshold rather than a global rejection threshold is chosen because under

1 the latter criteria strong deviations at certain locations (typically in the vicinity of pumping
2 wells) may be offset by small deviations at other wells.

3

4 For the retained simulators, the performance is assessed using equations (1), (2) and (3).

5 Point likelihood values estimated using equation (3) are combined into a global likelihood
6 value using a geometric mean inference function (Jensen, 2003). For each conceptual model,
7 predictive distributions of the state variables are obtained from the ensemble of likelihood
8 weighted predictions (rescaled such that for each conceptual model the likelihoods sum up to
9 one). Sampling from the prior input and parameter space continued until the first and second
10 moment of these predictive distributions stabilized.

11

12 **5. Results and discussion**

13 Since it is impossible to show the complete set of results for all combinations of likelihood
14 functions, variables, head observations, groundwater budget terms and alternative
15 conceptualizations, in the following sections the most relevant results are summarized.

16

17 **5.1. Convergence**

18 For the conceptual models 1Lhtg-L1 and 1Lhtg-L2 none of the simulations were accepted, as
19 all of them failed to meet the criteria of a maximum allowable departure of 5 m from the
20 observed heads. Hence, no results are presented for these models since they are discarded
21 from the posterior analysis.

22

23 Figure 2 shows, for the analysis based on the Gaussian likelihood function, the convergence
24 of the mean of the predictive distribution for the following groundwater budget terms: west
25 boundary condition (WBC) inflows (Figure 2a), recharge inflows (Figure 2b), west boundary
26 condition (WBC) outflows (Figure 2c), river gains (Figure 2d), and evapotranspiration (EVT)
27 outflows (Figure 2e). For all variables, convergence of the first moment was achieved in less

1 than 10,000 retained simulations, whereby groundwater budget terms converged to different
2 values in function of the conceptual model. Convergence of the second moment of the
3 predictive distributions, although not shown here, was also achieved within less than 10,000
4 retained simulations. It is important to note that the second moments converged to smaller
5 values when the alternative conceptual model approaches the true three-dimensional
6 hypothetical setup. These findings support the idea that predicting state variables relying on a
7 single conceptual model is prone to statistical bias and may produce an overconfident
8 estimation of predictive uncertainty. Similar patterns of convergence of the first and second
9 moments were observed for the other likelihood functions.

10

11 **5.2. Likelihood response surfaces**

12 Figure 3 shows the global likelihood response surface projected in one dimension for the six
13 unknown variables (see Table 2). The vertical dashed lines represent the true values used in
14 the three-dimensional hypothetical setup. It is seen from this figure that model performance is
15 highly sensitive to variables RECH and CH, as expressed by the well defined regions of
16 attraction centered on the true values (Figures 3a and 3b). For the other variables, well
17 performing simulators are found across the whole prior space. However, for EVTR and
18 SURF, zones of higher attraction are distinguished near their respective true values.

19

20 Figure 4 shows the normalized global likelihood response surface projected in two
21 dimensions for different combinations of normalized variables for model conceptualizations
22 1Lhtg-AVG, 2Lhtg, and 3Lhtg. For each combination the highest five normalized likelihood
23 values, which nearly all have a normalized likelihood larger than 0.95, are indicated by the
24 numbered white crosses. The two-dimensional projections reveal the complex nature of the
25 global likelihood response surface, with multiple localized zones of attraction and maximum
26 likelihood values located in different regions of the joint input and parameter space. This
27 reaffirms the idea of equifinality, i.e., that there exist multiple acceptable or behavioral

1 simulators that perform equally well and that can be spread over large regions of the model,
2 input and parameter space.

3

4 The plots in Figure 4 show that for increasing model complexity the regions of attractions
5 become more pronounced, or less diffuse, especially for the parameters to which model
6 performance is most sensitive. Plates a, b and c indicate that the two most sensitive
7 parameters (RECH and CH) are inversely correlated, with a tendency to a more defined
8 relationship with increasing model complexity. The other parameters do not show any strong
9 correlation, mainly due to the low sensitivity of the model performance to these parameters.

10

11 Results for the other likelihood functions, although not shown here, are very similar to those
12 presented in Figures 3 and 4. Hence, for the problem at hand, it can be concluded that the
13 shape of the likelihood response surface does not depend on the choice of likelihood function.

14

15 Figure 5 shows a one-dimensional projection of the likelihood response surface against the
16 model output variable river gains for three alternative conceptualizations (1Lhtg-AVG, 2Lhtg
17 and 3Lhtg) for the three likelihood functions used. It represents the weights (y-axis) that are
18 given to the different simulated values of river gain (x-axis) in the ensemble predictive
19 distribution of each model. As stated before, it is clear that the choice of likelihood does not
20 significantly impact the results. Increasing model complexity, on the other hand, results in a
21 slight increase of the maximum likelihood values, reduces the diffusivity of the likelihood
22 response surface and, for most groundwater budget terms, results in a more correct estimate
23 of the true values. Hence, although simpler models may result in simulations that are nearly
24 as good as the more complex models in terms of reproducing the set of head observations in
25 the training period, they typically lead to more bias and a larger predictive spread.

26

27 **5.3. Posterior model probabilities**

1 The integrated likelihoods and the posterior model probabilities of each alternative
2 conceptual model, approximated using equations (9) and (10), are presented in Table 3. The
3 posterior model probabilities represent the ability of each of the alternative models to
4 reproduce the observed data in the training period.

5
6 As previously stated, models 1Lhtg-L1 and 1Lhtg-L2 produced no results as none of the
7 simulations were able to meet the acceptance criteria. Hence, their integrated model
8 likelihood was set to zero and they were discarded in the calculation of the model ensemble
9 predictive distribution. In Table 3 it is seen that posterior model probability of the other
10 alternative conceptual models increases slightly from 0.18 to 0.22 with increasing level of
11 model complexity. The small difference in posterior model probability implies that, for the
12 given setup, the head observations do not allow to make a further distinction in performance
13 between the five retained conceptualizations. These results confirm that in real applications,
14 where the true hydrological concept is unknown and conditioning data are typically limited to
15 (a sparse set of) head observations, confining the model space to a single model is often not
16 supported by the data, hence advocate the idea of considering multiple conceptualizations. To
17 overcome this problem, other sources of qualitative or quantitative conditioning data that
18 allow a further differentiation between the models may be considered.

19

20 **5.4. Predictive distributions**

21 The posterior model probabilities are then used to combine the predictive distributions of the
22 five retained conceptual models using equation (4). The moments of the multi-model
23 ensemble predictive distribution are obtained through equations (7) and (8).

24

25 The cumulative predictive distributions of the groundwater budget terms for the five retained
26 conceptual models and the Bayesian model averaging are presented in Figure 6 for the
27 analysis based on the Gaussian likelihood function. The vertical dashed lines indicate the true

1 values observed from the three-dimensional hypothetical setup. Summarizing statistics of the
2 respective predictive distributions in Figure 6 are presented in Figure 7. Shown here for the
3 groundwater budget terms are the min, max and median values, as well as the inter-quartile
4 confidence intervals. Here the true values are represented by the horizontal lines.

5
6 Results for the individual models show that the predictive distributions of the budget terms
7 vary substantially in shape, central moment and spread between the different
8 conceptualizations. In general, it is seen that when the alternative conceptual model
9 approaches the true three-dimensional hypothetical setup, confidence intervals become
10 smaller and predictions are less biased, i.e., the median of the predictive distribution more
11 closely reproduces the observed value of the budget terms. These results show that although
12 the posterior model probabilities of the retained models differ only marginally their
13 predictions can vary substantially. Whereas for some of the (mainly simpler) models the true
14 groundwater budget values are not contained by the inter-quartile ranges of the predictions,
15 they are always captured by the inter-quartile range of the BMA ensemble predictions. This
16 reaffirms that relying on a single conceptual model is prone to statistical bias and may
17 produce an overconfident estimation of predictive uncertainty. The BMA on the other hand
18 provides consensus predictions and yields a more reliable estimation of the predictive
19 uncertainty.

20
21 The contribution of model uncertainty to the total predictive uncertainty is estimated using
22 equation (8) and is presented in Figure 8. Here, the total predictive variance is divided in
23 within-model and between-model variance for the five groundwater budget terms. Both
24 components are expressed as a percentage of the total variance. It is seen from this figure that
25 predictive variance due to the uncertainty in the conceptual model (between-model) ranges
26 from 5% for WBC outflows to approximately 30% for river gains, with practically no
27 difference between the results for the different likelihood functions. Information about the

1 susceptibility to conceptual model uncertainty of the different groundwater budget terms
2 provides useful information for possible improvement of the model concept or to guide
3 further data collection campaigns to optimally reduce conceptual uncertainty.

4

5 **6. Conclusions**

6 We presented a methodology to assess uncertainty in predictions of groundwater models
7 arising from errors in the model structure, forcing data and parameter estimates. The
8 methodology is based on the concept that there exist many good simulators of the system that
9 may be located in different regions of the combined model, input and parameter space, given
10 the data at hand. For a set of plausible system conceptualizations, input and parameter
11 realizations are sampled from the joint prior input and parameter space. A likelihood measure
12 is then calculated for each simulator based on its ability to reproduce system state variable
13 observations. The integrated likelihood of each conceptual model is obtained by integration
14 over the input and parameter space the likelihood of the different simulators. The integrated
15 likelihoods are consequently used in Bayesian model averaging to weight the model
16 predictions to obtain ensemble predictions.

17

18 The adopted approach is flexible in the sense that (i) there is no limitation in the number or
19 complexity of conceptual models that can be included, or to what degree input and parameter
20 uncertainty can be incorporated, (ii) any quantitative or qualitative (e.g., pumping well never
21 dries out) information about the system can be used to distinguish between different
22 simulators, (iii) the closeness between the predictions and system observations can be defined
23 in a variety of ways, including a formal statistical measure, and (iv) likelihoods, model
24 probabilities and predictive distributions can be easily updated when new information
25 becomes available. The major drawback of the approach is the computational burden inherent
26 to any Monte Carlo method.

1 For illustrative purposes the methodology was applied to a three-dimensional hypothetical
2 setup consisting of two aquifers separated by an aquitard, in which the flow field was
3 considerably affected by pumping wells and spatially variable hydraulic conductivity. A set
4 of 16 head observations sampled from this setup was used as conditioning data. The
5 proximity of the simulations to these observations was evaluated using three different
6 likelihood functions, including a formal statistical one. Seven alternative conceptualizations
7 with increasing complexity were adopted and only uncertainty in parameters and inputs that
8 were common to all conceptual models were considered. Two of the simpler one-layer
9 models were discarded from the further analysis as they failed to meet a subjectively chosen
10 criterion of closeness between the simulated and observed heads. For the other
11 conceptualizations convergence of the first and second moment of the predicted variable
12 distributions was achieved in less than 10,000 retained simulations.

13

14 The global likelihood response surface showed to be very complex, with multiple regions of
15 high likelihood and local maxima in different regions of the joint model, input and parameter
16 space. This confirms the concept of equifinality, i.e., that there exist many acceptable system
17 representations that cannot be easily rejected and that should be considered in assessing the
18 uncertainty associated with predictions. The likelihood response surfaces showed very little
19 dependence on the choice of the likelihood function adopted. As such, the selection of the
20 likelihood function did not have a significant impact on the further analysis and the general
21 patterns observed in the results were identical for the three likelihood functions.

22

23 The integrated likelihoods of the five retained models increased slightly with increasing
24 model complexity. The small differences in posterior model probability indicate that the set
25 of 16 head observations did not allow a further discrimination between the five retained
26 models. Nevertheless, predictive distributions of groundwater budget terms showed to be
27 considerably different in shape, central moment and spread among the models. When the

1 alternative conceptual model approached the true three-dimensional hypothetical setup,
2 confidence intervals were in general smaller and predictions were less biased. BMA, on the
3 other hand, provided consensus predictions yielding a more reliable estimation of the
4 predictive uncertainty. The contribution of model uncertainty to the total predictive
5 uncertainty varied between 5 to 30% depending on the groundwater budget term. The relative
6 contribution of model uncertainty for the different groundwater budget terms provides useful
7 information for updating the model concept or guiding data collection to optimally reduce
8 conceptual uncertainty.

9

10 The results of this study strongly advocate the idea to address conceptual model uncertainty
11 in the practice of groundwater modeling. With a hypothetical example it was shown that a set
12 of head observations, which in reality may often be the only information available about the
13 system dynamics, did not allow discriminating between a set of five models ranging from a
14 simple one-layer model to a conceptualization approaching the true three-dimensional setup.
15 Nevertheless, predictions of groundwater budget terms differed considerably among these
16 models. The use of a single model may result in smaller uncertainty intervals, hence an
17 increased confidence in the model simulations, but is very likely prone to statistical bias.
18 Also, in the presence of conceptual model uncertainty, which per definition can not be
19 excluded, this gain in accuracy in the short-term may have serious implications when using
20 the model for long-term predictions in which the system is subject to new stresses. It is
21 therefore advisable to explore a number of alternative conceptual models to obtain consensus
22 predictions that are more conservative, hence that are more likely to bracket the true system
23 responses.

24

25 It is expected that including other qualitative or quantitative sources of conditioning data,
26 such as conductivity data, geological profiles, transient groundwater head information, or

1 recharge estimates will allow a better differentiation between alternative models to further
2 reduce model uncertainty. These topics will be subject of future research.

3

4 **Acknowledgments**

5 The first author wishes to thank the Katholieke Universiteit (K.U.) Leuven for providing
6 financial support in the framework of PhD IRO-scholarships. We also wish to acknowledge
7 the assistance provided by Jan de Laet, and Wim Obbels with running codes on the K.U.
8 Leuven supercomputer (VIC CLUSTER) and by Jorge Gonzalez to implement the R scripts.
9 Comments and suggestions by Marijke Huysmans and Okke Batelaan for improving the
10 original manuscript are greatly appreciated.

1 **References**

- 2 Abraham, R., and L. See (2002), Multi-model data fusion for river flow forecasting: an
3 evaluation of six alternative methods based on two contrasting catchments, *Hydrology and*
4 *Earth System Sciences*, 6(4), 655-670.
- 5
- 6 Ajami, N., Q. Duan, X. Gao and S. Sorooshian (2005), Multi-model combination techniques
7 for hydrologic forecasting: application to distributed model intercomparison project results,
8 *Journal of Hydrometeorology*, 7(4), 755-768.
- 9
- 10 Akaike, H. (1974), A new look at the statistical model identification, *IEEE Transactions on*
11 *Automatic Control* AC 19, 716-723.
- 12
- 13 Bernardo, J. and A. Smith (2000), *Bayesian theory*, 608 pp., John Wiley & Sons Inc,
14 Chichester.
- 15
- 16 Beven, K. and A. Binley (1992), The future of distributed models – model calibration and
17 uncertainty prediction, *Hydrological Processes*, 6(3), 279-283.
- 18
- 19 Beven, K. (1993), Prophecy, reality and uncertainty in distributed hydrological modeling,
20 *Advances in Water Resources*, 16(1), 41-51.
- 21
- 22 Beven, K. and J. Freer (2001), Equifinality, data assimilation, and uncertainty estimation in
23 mechanistic modelling of complex environmental systems using the GLUE methodology,
24 *Journal of Hydrology*, 249(1-4), 11-29.
- 25
- 26 Beven, K. (2005a), On the concept of model structural error, *Water Science & Technology*,
27 52(6), 167-175.

1 Beven, K. (2005b), A manifesto for the equifinality thesis, *Journal of Hydrology*, 320(1-2),
2 18-36.
3
4 Beven, K., P. Smith and J. Freer (2007), Comment on “Hydrological forecasting uncertainty
5 assessment: Incoherence of the GLUE methodology” by Pietro Mantovan and Ezio Todino,
6 *Journal of Hydrology*, doi: 10.1016/j.jhydrol.2007.02.023.
7
8 Binley, A. and K. Beven (2003), Vadose zone flow model uncertainty as conditioned on
9 geophysical data, *Ground Water*, 41(2), 119-127.
10
11 Bredehoeft, J. (2003), From models performance assessment: The conceptualization problem,
12 *Ground Water*, 41(5), 571-577.
13
14 Bredehoeft, J. (2005), The conceptualization model problem – surprise, *Hydrogeology*
15 *Journal*, 13(1), 37-46.
16
17 Burnham, K. and D. Anderson (1998), *Model selection and inference. A practical*
18 *information-theoretic approach*, 353 pp., Springer-Verlag, New York.
19
20 Carrera, J., A. Alcolea, A. Medina, J. Hidalgo and L. Slooten (2005), Inverse problem in
21 hydrogeology, *Hydrogeology Journal*, 13(1), 206-222.
22
23 Deutsch, C. and A. Journel (1998), *GSLIB Geostatistical software library and user’s guide*,
24 2nd ed., Oxford University Press, New York.
25
26 Draper, D. (1995), Assessment and propagation of model uncertainty, *Journal of the Royal*
27 *Statistical Society Series B (with discussion)*, 57(1), 45-97.

1 Feyen, L., K. Beven, F. De Smedt and J. Freer (2001), Stochastic capture zone delineation
2 within the GLUE-methodology: conditioning on head observations, *Water Resources*
3 *Research*, 37(3), 625-638.
4
5 Feyen, L., P. Ribeiro, F. De Smedt, P. Diggle (2002), Bayesian methodology to stochastic
6 capture zone determination: conditioning on transmissivity measurements, *Water Resources*
7 *Research*, 38(9), doi:10.1029/2001WR000950.
8
9 Freer, J. and K. Beven (1996), Bayesian estimation of uncertainty in runoff prediction and the
10 value of data: An application of the GLUE approach, *Water Resources Research*, 32(7),
11 2161-2173.
12
13 Georgakakos, K., D. Seo, H. Gupta, J. Schaake and M. Butts (2004), Characterizing
14 streamflow simulation uncertainty through multimodel ensembles, *Journal of Hydrology*,
15 298(1-4), 222-241.
16
17 Harbaugh, A., E. Banta, M. Hill and M. McDonald (2000), MODFLOW-2000, U.S.
18 Geological Survey modular ground-water model-user guide to modularization concepts and
19 the ground-water flow process, U.S. Geol. Surv. *Open File Rep.*, 00-92, 121 pp.
20
21 Hoeting, J., D. Madigan, A. Raftery and C. Volinsky (1999), Bayesian model averaging: A
22 tutorial, *Statistical Science*, 14(4), 382-417.
23
24 Jensen, J. (2003), Parameter and uncertainty estimation in groundwater modelling, PhD
25 thesis, 143 pp., Aalborg University, Aalborg, Dec 2002.
26

1 Kashyap, R. (1982), Optimal choice of AR and MA parts in autoregressive moving average
2 models, IEEE Transactions Pattern Analysis and Machine Intelligence, PAMI 42(9), 99-104.
3

4 Madigan, D. and A. Raftery (1994), Model selection and accounting for model uncertainty in
5 graphical models using Occam's window, Journal of the American Statistical Association,
6 89(428), 1535-1546.
7

8 Mantovan, P. and E. Todini (2006) Hydrological forecasting uncertainty assessment:
9 Incoherence of the GLUE methodology, Journal of Hydrology, 330(1-2), 368-381.
10

11 Morse, G., G. Pohll, J. Huntington and R. Rodriguez (2003), Stochastic capture zone analysis
12 of an arsenic-contaminated well using the generalized likelihood uncertainty estimator
13 (GLUE) methodology, Water Resources Research, 39(6), 1151, doi:10.1029/2002WR001470
14

15 Nash, J. and J. Sutcliffe (1970) River flow forecasting through conceptual models. Part I – A
16 discussion of principles, Journal of Hydrology, 10(3), 282-290.
17

18 Neuman, S. (2003), Maximum likelihood Bayesian averaging of uncertain model predictions,
19 Stochastic Environmental Research and Risk Assessment, 17(5), 291-305.
20

21 Neuman, S. and P. Wierenga (2003), A comprehensive strategy of hydrogeologic modelling
22 and uncertainty analysis for nuclear facilities and sites, Report NUREG/CR-6805, US
23 Nuclear Regulatory Commission, Washington, USA.
24

25 Poeter, E. and D. Anderson (2005), Multimodel ranking and inference in ground water
26 modeling, Ground Water, 43(4), 597-605.
27

1 Raftery, A. and Y. Zheng (2003), Discussion: Performance of Bayesian model averaging,
2 Journal of the American Statistical Association, 98(464), 931-938.
3
4 Refsgaard, J., J. van der Sluijs, J. Brown and P. van de Keur (2006), A framework for dealing
5 with uncertainty due to model structure error, Advances in Water Resources, 29(11), 1586-
6 1597.
7
8 Romanowicz, R., K. Beven and J. Tawn (1994), Evaluation of prediction uncertainty in non-
9 linear hydrological models using a Bayesian approach, in *Statistics for the Environment II;*
10 *Water Related Issues*, edited by V. Barnett and K. Turkman, pp. 297-317, Wiley, New York.
11
12 Rubin, Y. (2003), *Applied stochastic hydrogeology*, 391 pp., Oxford University Press, New
13 York.
14
15 Sun, N-Z. (1994), *Inverse problems in groundwater modelling. Theory and applications of*
16 *transport in porous media*, 352 pp., Kluwer Academic Publishers, Netherlands.
17
18 Ye, M., S. Neuman and P. Meyer (2004), Maximum likelihood Bayesian averaging of spatial
19 variability models in unsaturated fractured tuff, Water Resources Research, 40, W05113,
20 doi:10.1029/2003WR002557.
21
22 Ye, M., S. Neuman, P. Meyer and K. Pohlmann (2005), Sensitivity analysis and assessment
23 of prior model probabilities in MLBMA with application to unsaturated fractured tuff, Water
24 Resources Research, 41, W12429, doi:10.1029/2005WR004260.

1 **Figure captions**

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

4

5 Figure 2: Convergence of the first moment of the predictive distributions of the groundwater
6 budget terms as a function of the Number of retained Monte Carlo Simulations (NMCS) for
7 the Gaussian likelihood function (GAUSS): (a) west boundary condition (WBC) inflows, (b)
8 recharge inflows, (c) west boundary condition (WBC) outflows, (d) river gains and (e)
9 evapotranspiration (EVT) outflows.

10

11 Figure 3: One-dimensional projection of the global likelihood response surface (based on the
12 Gaussian likelihood function) for the six parameters for conceptual model 3Lhtg. Vertical
13 dashed lines represent the parameter values used in the three-dimensional hypothetical setup.

14

15 Figure 4: Two-dimensional projection of the normalized likelihood response surface (based
16 on the Gaussian likelihood function) for the normalized parameters RECH vs. CH, RECH vs.
17 EVTR and, RECH vs. RIVC for the alternative conceptual models 1Lhtg-AVG, 2Lhtg and
18 3Lhtg. Numbered crosses represent the locations of the five highest likelihood values.

19

20 Figure 5: Results for the river gains for the alternative conceptual models 1Lhtg-AVG (a-d-
21 g), 2Lhtg (b-e-h) and 3Lhtg (c-f-i), and the Gaussian – GAUSS (a-c), Triangular – TRIANG
22 (d-e) and Model efficiency – MODEFF (g-i) likelihood functions. Vertical dashed-lines
23 represent the observed values from the three-dimensional hypothetical setup.

24

25 Figure 6: Cumulative probability distributions of the groundwater budget terms for the five
26 alternative conceptual models and the Bayesian model averaging (BMA) based on the

1 Gaussian likelihood function: (a) west boundary condition (WBC) inflows, (b) recharge
2 inflows, (c) west boundary condition (WBC) outflows, (d) river gains and (e)
3 evapotranspiration (EVT) outflows. Vertical dashed-lines represent observed values from the
4 three-dimensional hypothetical setup.

5

6 Figure 7: Total variance estimated using equation (8) for the groundwater budget terms based
7 on the Gaussian (GAUSS), triangular (TRIANG) and model efficiency (MODEFF)
8 likelihood function. From left to right: west boundary condition (WBC) inflows, recharge
9 inflows, west boundary condition (WBC) outflows, river gains and evapotranspiration (EVT)
10 outflows.

11

12 Figure 8: Summary statistics of the predictive distributions of the alternative conceptual
13 models and multi-model BMA prediction for the groundwater budget terms: a) west
14 boundary condition (WBC) inflows, (b) recharge inflows, (c) west boundary condition
15 (WBC) outflows, (d) river gains and (e) evapotranspiration (EVT) outflows. Horizontal lines
16 represent the values obtained from the three-dimensional hypothetical setup. Q_1 and Q_3
17 represent the first and third quartile, respectively.

1 **Tables**

2 Table 1: Parameters describing the hydraulic conductivity spatial correlation structure for the
 3 different layers of the three-dimensional hypothetical setup (Based on Rubin (2003), Tables
 4 2.1 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

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

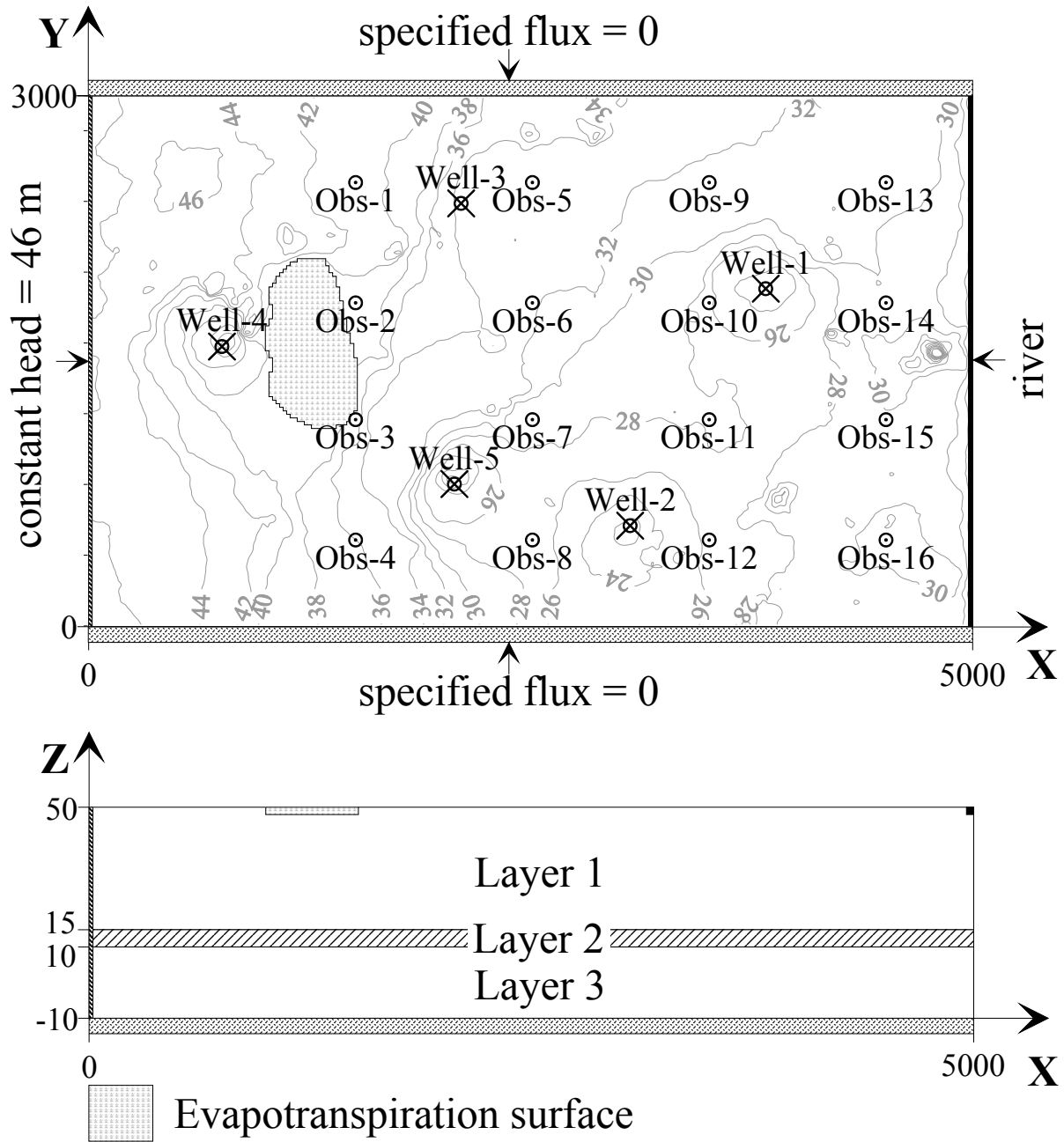
8

9

10 Table 3: Summary of the integrated likelihood and posterior model probabilities for the
 11 alternative conceptual models

	Likelihood function	Conceptual model							Total
		1Lhtg-L1	1Lhtg-L2	1Lhtg-L3	1Lhtg-AVG	2Lhtg	2LQ3Dhtg	3Lhtg	
$p(\mathbf{D} \mathbf{M}_k)$	GAUSS	0	0	902.6	935.6	990.4	1046.9	1079.4	4954.9
	TRIANG	0	0	4385.5	4608.1	4997.4	5365.2	5407.3	24763.5
	MODEFF	0	0	5952.6	6191.6	6579.3	6944.5	6994.7	32662.6
$p(\mathbf{M}_k)$		1/7	1/7	1/7	1/7	1/7	1/7	1/7	1.0
$p(\mathbf{M}_k \mathbf{D})$	GAUSS	0	0	0.1822	0.1888	0.1999	0.2113	0.2178	1.0
	TRIANG	0	0	0.1771	0.1861	0.2018	0.2167	0.2184	1.0
	MODEFF	0	0	0.1822	0.1896	0.2014	0.2126	0.2141	1.0

1 Figures

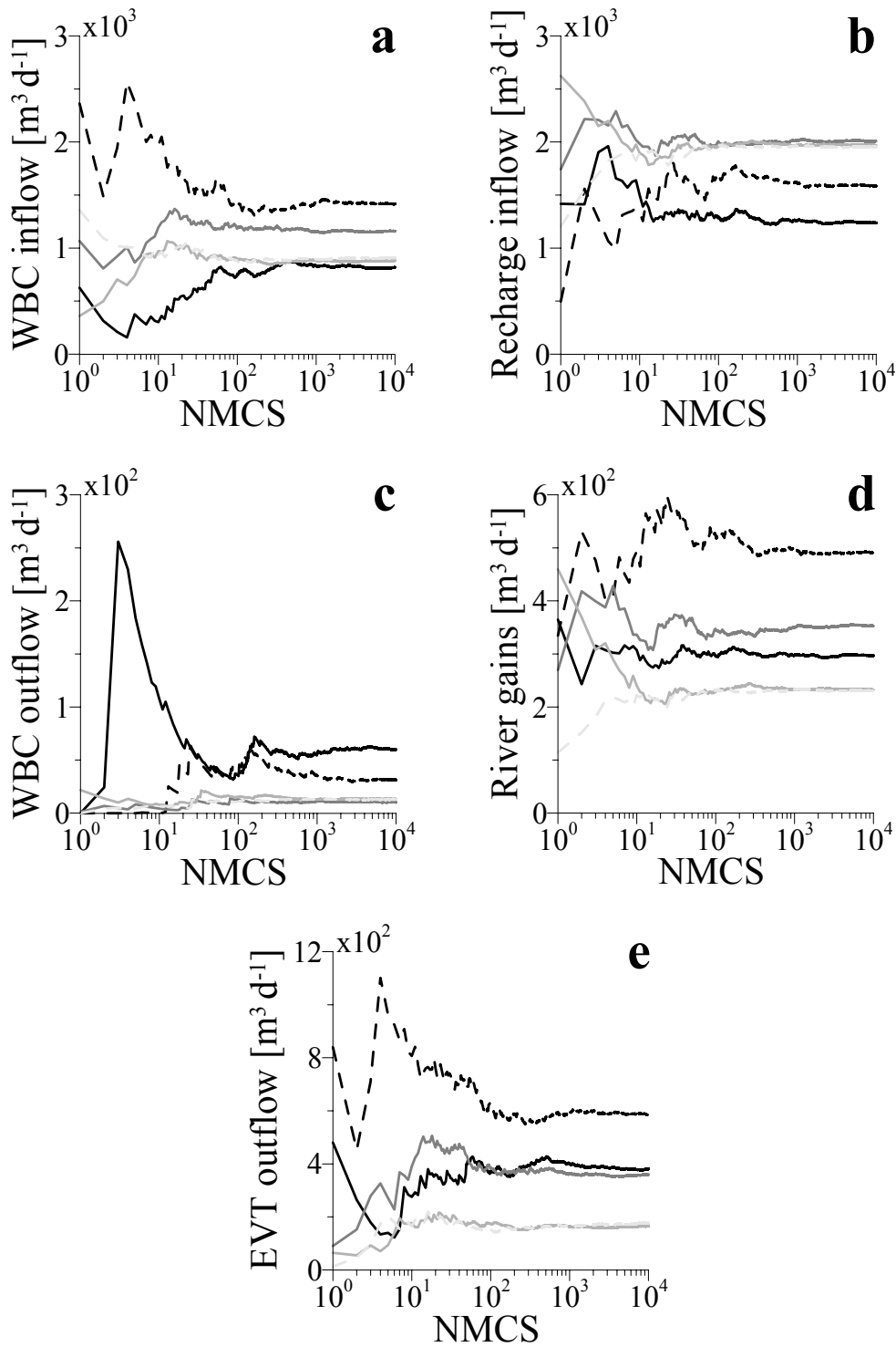


2

Evapotranspiration surface

3 Figure 1: Three-dimensional hypothetical setup including (O) observation wells and (X)

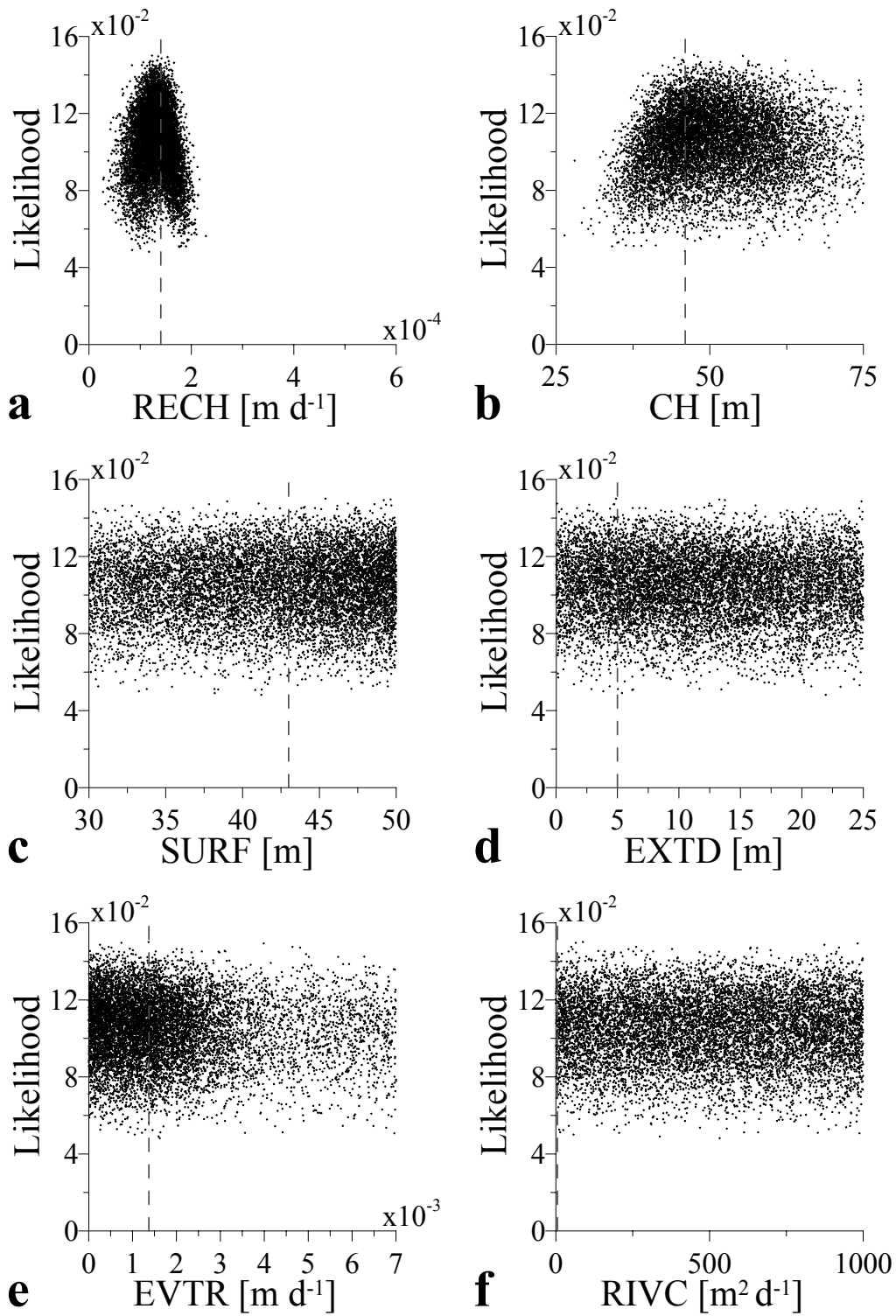
4 pumping wells overlain by the groundwater head distribution in the first layer.



— 1Lhtg-AVG - - 1Lhtg-L3 — 2Lhtg — 2LQ3Dhtg - - 3Lhtg

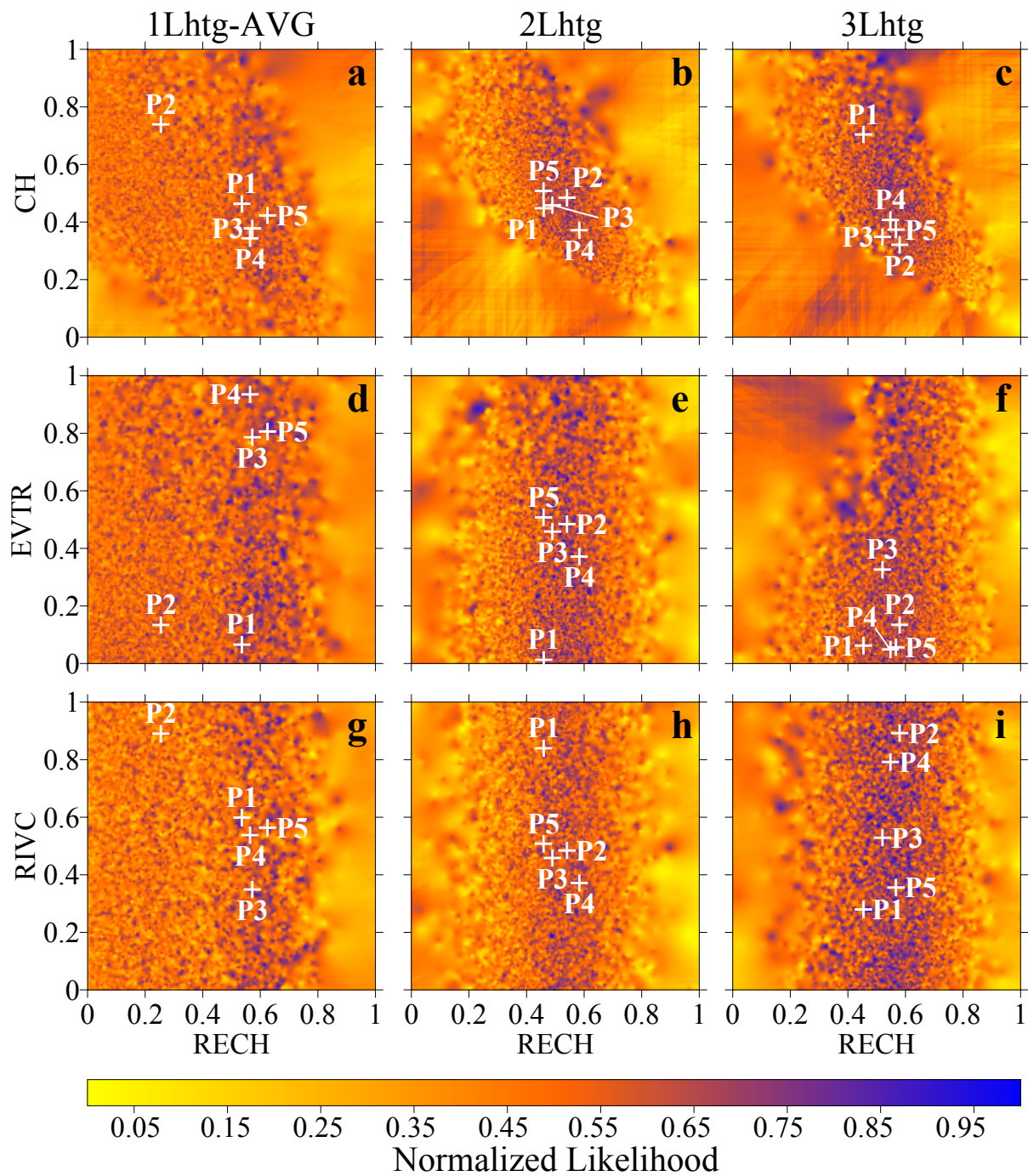
1

2 Figure 2: Convergence of the first moment of the predictive distributions of the groundwater
 3 budget terms as a function of the Number of retained Monte Carlo Simulations (NMCS) for
 4 the Gaussian likelihood function (GAUSS): (a) west boundary condition (WBC) inflows, (b)
 5 recharge inflows, (c) west boundary condition (WBC) outflows, (d) river gains and (e)
 6 evapotranspiration (EVT) outflows.

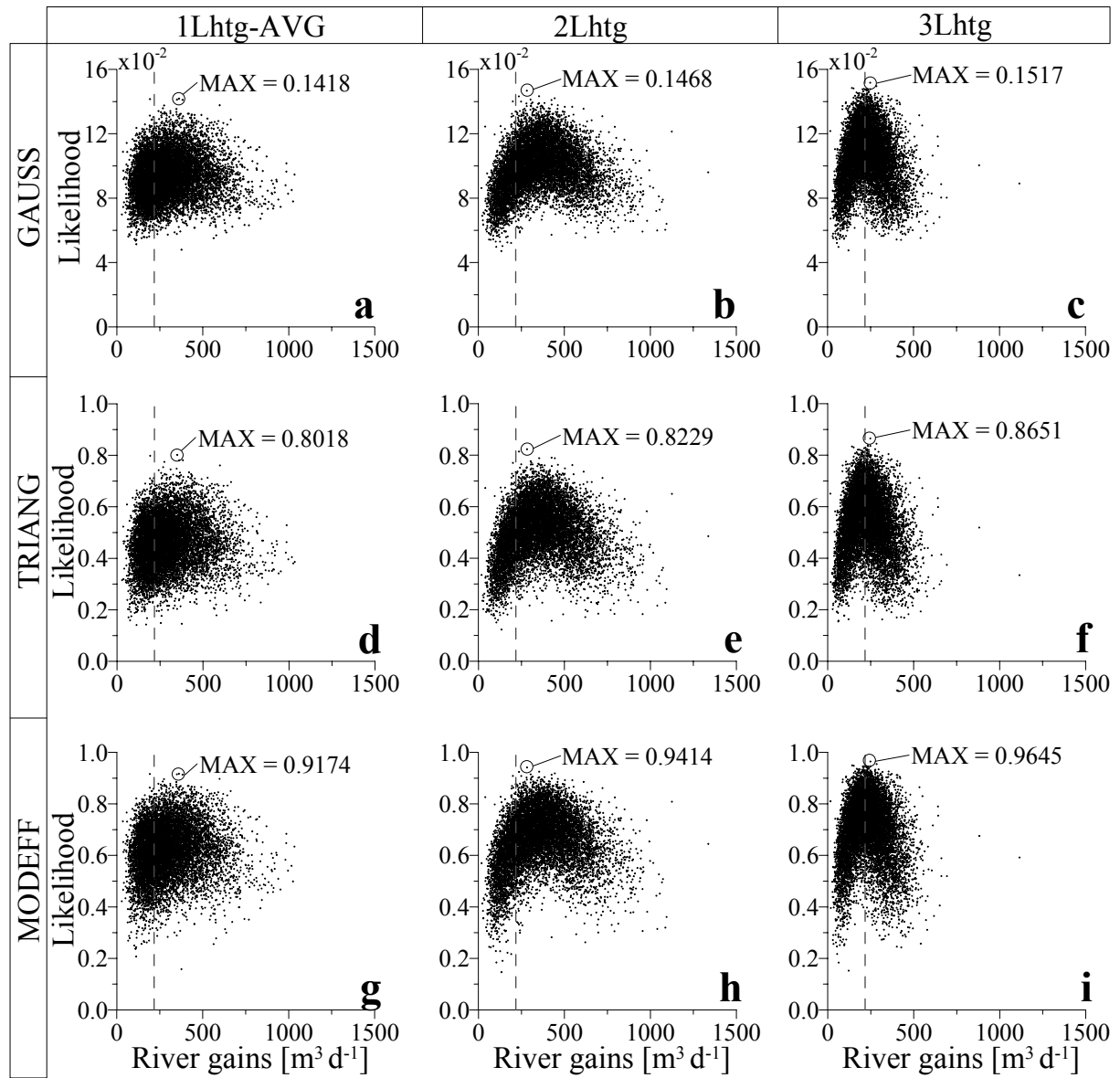


1

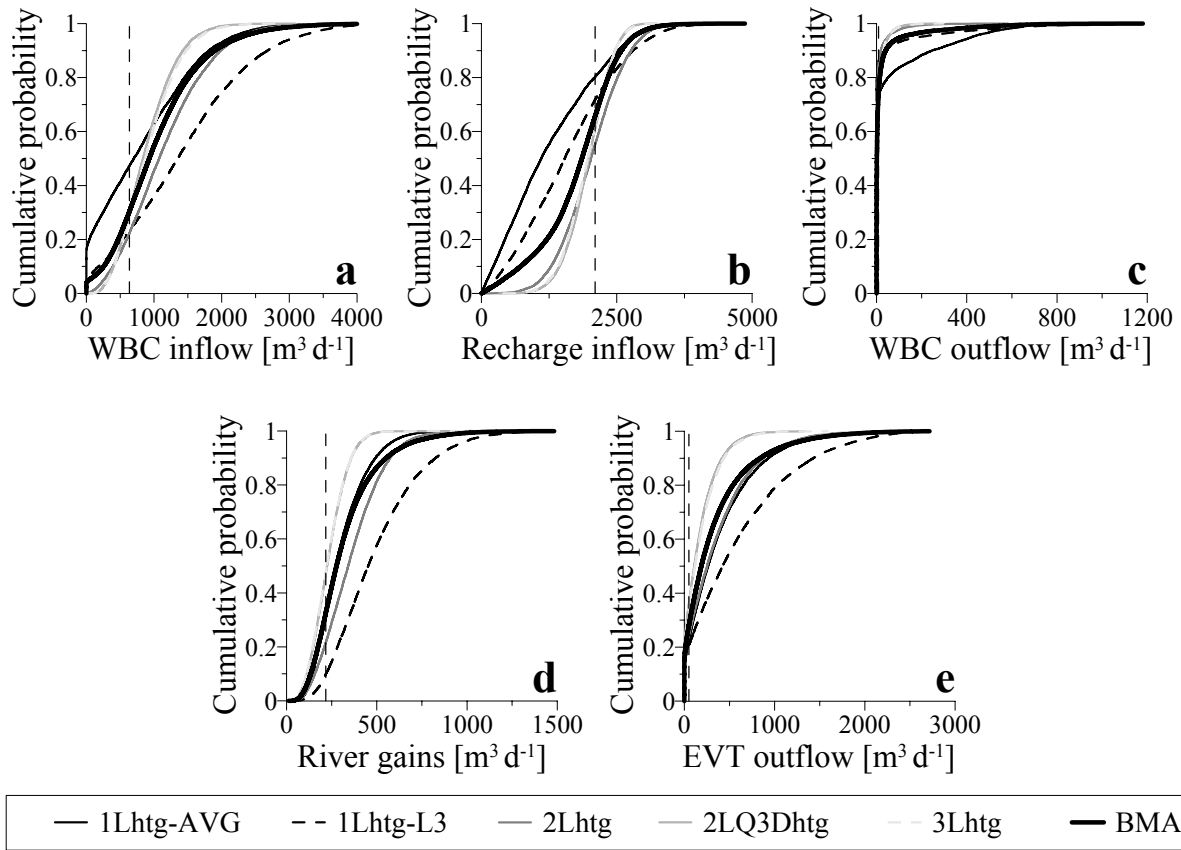
2 Figure 3: One-dimensional projection of the global likelihood response surface (based on the
 3 Gaussian likelihood function) for the six parameters for conceptual model 3Lhtg. Vertical
 4 dashed lines represent the parameter values used in the three-dimensional hypothetical setup.



1
 2 Figure 4: Two-dimensional projection of the normalized likelihood response surface (based
 3 on the Gaussian likelihood function) for the normalized parameters RECH vs. CH, RECH vs.
 4 EVTR and, RECH vs. RIVC for the alternative conceptual models 1Lhtg-AVG, 2Lhtg and
 5 3Lhtg. Numbered crosses represent the locations of the five highest likelihood values.



1
2 Figure 5: Results for the river gains for the alternative conceptual models 1Lhtg-AVG (a-d-
3 g), 2Lhtg (b-e-h) and 3Lhtg (c-f-i), and the Gaussian – GAUSS (a-c), Triangular – TRIANG
4 (d-e) and Model efficiency – MODEFF (g-i) likelihood functions. Vertical dashed-lines
5 represent the observed values from the three-dimensional hypothetical setup.



1 — 1Lhtg-AVG -- 1Lhtg-L3 — 2Lhtg — 2LQ3Dhtg - - 3Lhtg — BMA

2 Figure 6: Cumulative probability distributions of the groundwater budget terms for the five

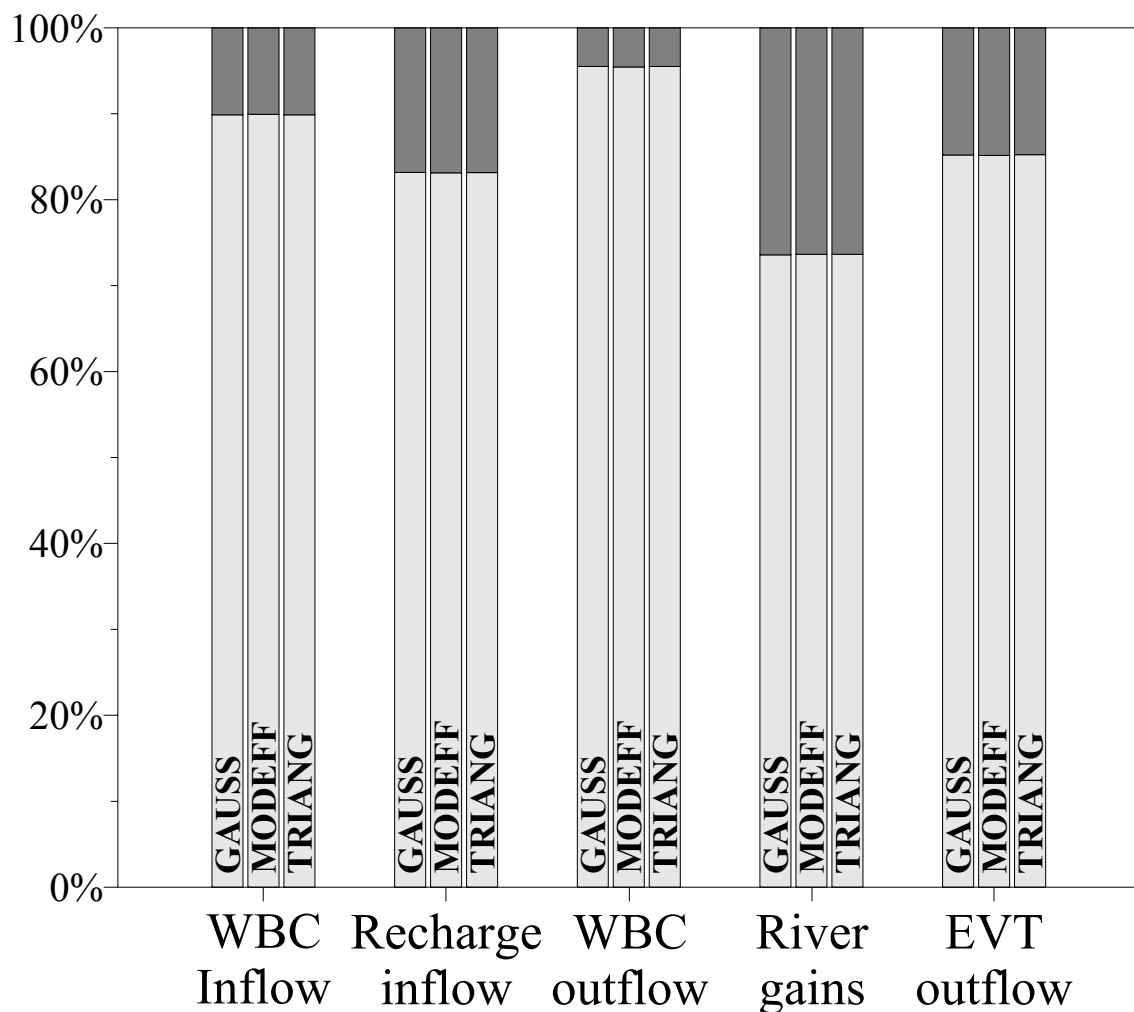
3 alternative conceptual models and the Bayesian model averaging (BMA) based on the

4 Gaussian likelihood function: (a) west boundary condition (WBC) inflows, (b) recharge

5 inflows, (c) west boundary condition (WBC) outflows, (d) river gains and (e)

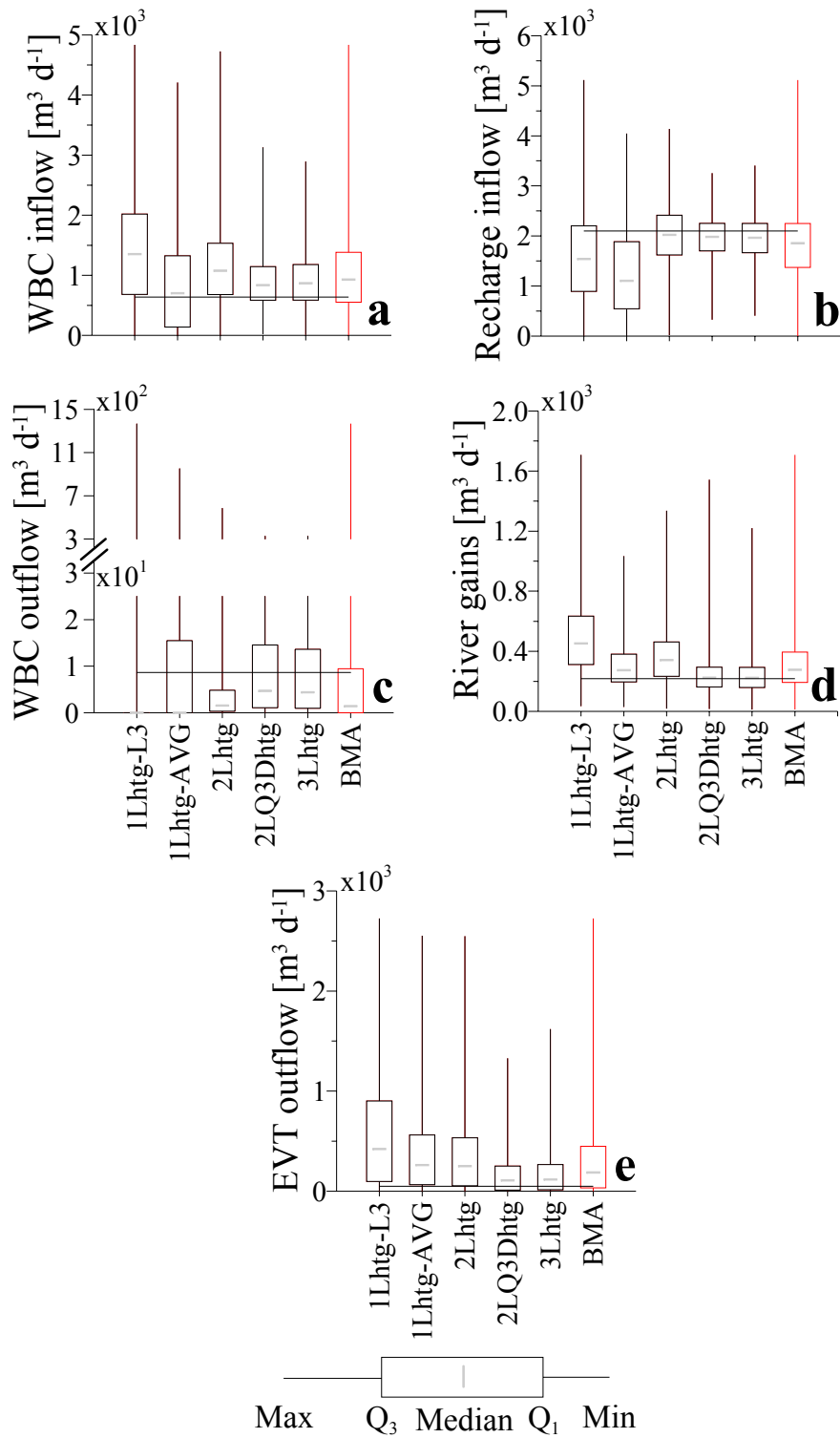
6 evapotranspiration (EVT) outflows. Vertical dashed-lines represent observed values from the

7 three-dimensional hypothetical setup.



Within-model variance
 Between-model variance

1
 2 Figure 7: Total variance estimated using equation (8) for the groundwater budget terms based
 3 on the Gaussian (GAUSS), triangular (TRIANG) and model efficiency (MODEFF)
 4 likelihood function. From left to right: west boundary condition (WBC) inflows, recharge
 5 inflows, west boundary condition (WBC) outflows, river gains and evapotranspiration (EVT)
 6 outflows.



1
 2 Figure 8: Summary statistics of the predictive distributions of the alternative conceptual
 3 models and multi-model BMA prediction for the groundwater budget terms: a) west
 4 boundary condition (WBC) inflows, (b) recharge inflows, (c) west boundary condition
 5 (WBC) outflows, (d) river gains and (e) evapotranspiration (EVT) outflows. Horizontal lines
 6 represent the values obtained from the three-dimensional hypothetical setup. Q₁ and Q₃
 7 represent the first and third quartile, respectively.