

1 **Assessing the effects of spatial discretization on large-scale flow model performance**
2 **and prediction uncertainty**

3 Wildemeersch, S^{a,*}, Goderniaux P.^b, Orban Ph.^a, Brouyère S.^a, Dassargues A.^a

4 ^a *University of Liege, ArGEnCo, GEO³, Hydrogeology and Environmental Geology, Aquapôle, B52/3 Sart-*
5 *Tilman, 4000 Liege, Belgium*

6 ^b *University of Mons, Faculty of Engineering, Fundamental and Applied Geology Department, Rue de*
7 *Houdain 9, 7000 Mons, Belgium*

8 * Corresponding author at: University of Liege, ArGEnCo, GEO³, Hydrogeology and Environmental
9 Geology, Aquapôle, B52/3 Sart-Tilman, 4000 Liege, Belgium. Tel.: +32 (0)43669553. Fax:
10 +32 (0)43669520.

11 *E-mail address:* swildemeersch@ulg.ac.be (S. Wildemeersch)

12 **Abstract**

13 Large-scale physically-based and spatially-distributed models (>100 km²) constitute useful tools for
14 water management since they take explicitly into account the heterogeneity and the physical processes
15 occurring in the subsurface for predicting the evolution of discharge and hydraulic heads for several
16 predictive scenarios. However, such models are characterized by lengthy execution times. Therefore,
17 modelers often coarsen spatial discretization of large-scale physically-based and spatially-distributed
18 models for reducing the number of unknowns and the execution times. This study investigates the
19 influence of such a coarsening of model grid on model performance and prediction uncertainty. The
20 improvement of model performance obtained with an automatic calibration process is also investigated.
21 The results obtained show that coarsening spatial discretization mainly influences the simulation of
22 discharge due to a poor representation of surface water network and a smoothing of surface slopes that

23 prevents from simulating properly surface water-groundwater interactions and runoff processes.
24 Parameter sensitivities are not significantly influenced by grid coarsening and calibration can
25 compensate, to some extent, for model errors induced by grid coarsening. The results also show that
26 coarsening spatial discretization mainly influences the uncertainty on discharge predictions. However,
27 model prediction uncertainties on discharge only increase significantly for very coarse spatial
28 discretizations.

29 *Keywords:* spatial discretization; model performance; sensitivity analysis; automatic calibration;
30 prediction uncertainty.

31 **1 Introduction**

32 Large-scale physically-based and spatially-distributed models ($> 100 \text{ km}^2$) are increasingly used in water
33 management for their unique capacity of gathering every piece of information obtained on a
34 hydrological system to simulate its quantitative and qualitative evolution for several predictive
35 scenarios. These models are intended to provide predictions on both the integrated response
36 (discharge) and the distributed response (hydraulic heads) of the catchment.

37 Physically-based and spatially-distributed models take explicitly into account the heterogeneity and the
38 physical processes occurring in the surface and the subsurface. Therefore, they are expected to provide
39 predictions with higher level of confidence than black-box models (e.g. Ebel and Loague, 2006; Li et al.;
40 2008; Goderniaux et al., 2009). Additionally, they are also used for improving the understanding of the
41 physics of hydrological processes (e.g. Frei et al., 2009; Meyerhoff and Maxwell, 2011; Irvine et al.,
42 2012). However, physically-based and spatially-distributed models are characterized by lengthy
43 execution times, especially for integrated surface and subsurface transient flow simulations at large-
44 scale. Consequently, choices and simplifications are made for obtaining tractable execution times. The

45 most common simplification consists in coarsening the spatial discretization for reducing the number of
46 unknowns of the problem and the execution time. The effects of such a coarsening of model grid are
47 worthy being studied since they can limit the accuracy of model results and increase model prediction
48 uncertainties.

49 A series of studies have already been performed on the effects of spatial discretization on physically-
50 based and spatially-distributed model performance. Refsgaard (1997) calibrated and validated a 3D
51 model with a 500 m grid i.e. with a constant element size of 500 m (no refinement) for the Karup
52 catchment in Denmark (440 km²). Three other models with 1000 m, 2000 m, and 4000 m grids were
53 then generated using the same parameter values than those obtained by calibration for the initial model
54 (no recalibration, no upscaling). The models were compared in terms of both discharge and hydraulic
55 heads. The results from this study indicated that runoff was poorly simulated by the models coarser
56 than 1000 m due to a poor representation of the surface water network which prevents from simulating
57 properly surface water-groundwater interaction. However, the author suggested that a significant
58 recalibration of models with a coarse grid could improve their performance. This is supported by the
59 study of Vázquez et al. (2002). They calibrated a 3D model with a 600 m grid i.e. with a constant element
60 size of 600 m (no refinement) for the Gete catchment in Belgium (586 km²). They also generated a
61 second model with a finer grid (300 m) and a third model with a coarser grid (1200 m) using the same
62 parameters than those obtained by calibration for the initial model (no recalibration, no upscaling).
63 These 300 m and 1200 m grid models were then recalibrated individually using a trial-and-error
64 calibration process. As for the study of Refsgaard (1997), the models were compared in terms of both
65 discharge and hydraulic heads. Although, in general, model results remained worse for the 1200 m grid
66 model than for the 300 m and the 600 m grid models, this study proved that a recalibration is required
67 for obtaining effective parameter values and improving model performance when the grid resolution is
68 changed. Sciuto and Diekkrüger (2010) developed a 3D model with a 25 m grid refined in the river zone

69 for the Wüstebach catchment in Germany (0.27 km²). They also developed a second model with a 100 m
70 grid using the mean averaging method for upscaling parameter values and a third model with the same
71 model grid than the initial model and the same soil configuration than the second model. They
72 compared the results obtained in terms of discharge and spatial pattern of soil moisture. The influence
73 of upscaling was investigated by comparing the first and the second models and the effects of spatial
74 discretization were studied by comparing the second and the third models. They showed that a coarse
75 grid leads to higher discharge and less actual evapotranspiration than a fine grid due to the smoothing
76 of soil surface which induces a loss of topographic information. They also showed that the upscaling
77 technique they selected was efficient for simulating discharge and spatial pattern of soil moisture. They
78 suggested that the nonlinear relationship between soil moisture and evapotranspiration could explain
79 the deterioration of model results when the grid is coarsened without parameter upscaling. However,
80 none of their models were calibrated. Downer and Ogden (2004) performed a spatial convergence study
81 for the Hortonian Godwin Creek Experimental catchment (21.2 km²) and the non-Hortonian Muddy
82 Brook catchment (3.64 km²) in the US. They developed a series of 2D vadose zone model of increasing
83 vertical cell size for each of these catchments. The models were calibrated with an automated
84 calibration process using the shuffled complex evolution method. The calibrated models were compared
85 in terms of infiltration, runoff, and evapotranspiration fluxes to evaluate the appropriate vertical
86 discretization required for accurately solving the Richards' equation. The results from this study showed
87 that small vertical cell size (on the order of centimetres) is required in the unsaturated zone to
88 accurately simulate hydrological fluxes. However, providing that effective parameters obtained by
89 calibration are used, the results of this study also shows that it is possible to slightly increase vertical cell
90 size in the unsaturated zone without significantly deteriorating the simulation of hydrological fluxes.
91 These results about the vertical cell size required in the unsaturated zone for accurately solving the
92 Richards' equation are consistent with those obtained by Vogel and Ippisch (2008).

93 All these studies provide valuable information on the effects of spatial discretization on model
94 performance. However, most of them neglect the calibration or use a simple trial-and-error calibration
95 process which is, by nature, subjective (Poeter and Hill, 1997). An automatic calibration process is
96 essential for properly evaluating the capacity of calibration to improve the performance of models with
97 a coarse grid. The present study includes such an automatic calibration process. Additionally, the
98 present study includes for the first time an evaluation of the influence of spatial discretization on model
99 prediction uncertainties by comparing the linear confidence intervals on predictions calculated for each
100 model.

101 The objective of the present study is to evaluate the effects of horizontal spatial discretization on
102 discharge and hydraulic heads simulated by a large-scale physically-based and spatially-distributed
103 model. This evaluation is performed using graphs of model fit and performance criteria. The
104 improvement of model performance obtained with an automatic calibration process is also investigated
105 and linear confidence intervals on predictions are calculated for each model. The results of this study
106 can help modelers defining the horizontal spatial discretization for their models by better perceiving its
107 influence on model performance and model prediction uncertainties.

108 **2 Methodology**

109 The effects of horizontal spatial discretization on model performance and model prediction
110 uncertainties are investigated using a synthetic catchment. The hydrological processes in this synthetic
111 catchment are simulated with HydroGeoSphere (Therrien et al., 2012). HydroGeoSphere is a fully-
112 integrated physically-based hydrological model capable of solving very complex problems such as
113 integrated flow in large-scale catchments (for example, see Goderniaux et al., 2009; 2011). Two-
114 dimensional surface water flow is represented using the two-dimensional diffusion-wave approximation
115 to the Saint-Venant equation. Three-dimensional subsurface water flow in both the saturated and the

116 vadose zones is represented using the Richards' equation. The processes of interception and
117 evapotranspiration are modeled following the conceptualization of Kristensen and Jensen (1975). The
118 coupling of the surface to the subsurface is either performed with the common node approach
119 (continuity of hydraulic head between the two domains) or the dual node approach (exchange of water
120 between the two domains via a first-order exchange coefficient). A complete description of
121 HydroGeoSphere is available in Therrien et al. (2012). A short summary is provided in the paper of Li et
122 al. (2008) and in the software spotlight of Brunner and Simmons (2012a).

123 The choice of working with a synthetic catchment instead of a real catchment is motivated by the wish
124 of focusing only on the effects of horizontal spatial discretization on model performance. When working
125 with a synthetic catchment, the model geometry, the parameter values, and the boundary conditions
126 are exactly known. Furthermore, there is no measurement error on the observations produced.
127 Therefore, it is possible to test specific model features such as the influence of grid resolution on
128 discharge and hydraulic head simulation without unintentionally taking into account other sources of
129 errors related to a lack of knowledge of the hydrological system. The concept of synthetic catchment is
130 quite usual in hydrogeology (for example, see Poeter and McKenna, 1995; Hill et al., 1998; Schäfer et al.,
131 2004; Bauer et al., 2006, Beyer et al., 2006). The synthetic catchment generated for this study is complex
132 in that the flow system is fully-integrated and physically-based with consistent physical state
133 parameters. However, the synthetic catchment is simplified with respect to the heterogeneity of land
134 use and geology in reality. Yet, this study focuses on the effects of spatial discretization on model
135 performance and not on the influence of heterogeneity representation. The way grid size influences
136 model results would have been similar for a synthetic catchment with a higher level of heterogeneity,
137 provided that the heterogeneity is correctly represented. Therefore, despite this simplification, the
138 synthetic catchment is judged complex enough to serve the objective of this study.

139 The methodology involves three main steps:

140 **STEP 1 – Generation of the reference model/Generation of models with a coarse grid.** A 5-year
141 simulation is run with the reference model for producing reference discharge and hydraulic head
142 observations. The reference model is characterized by a fine spatial discretization. The same 5-year
143 simulation is then run with models with a coarse grid using the same parameter values than those used
144 in the reference model (no calibration). These models with a coarse grid differ by their horizontal spatial
145 discretization (constant element size of 250 m, 500 m, 750 m model, or 1000 m). The simulated values
146 of discharge and hydraulic head obtained with these models are saved for further graphical model fit
147 analysis and calculation of performance criteria.

148 **STEP 2 – Calibration of models with a coarse grid.** The models with a coarse grid are individually
149 calibrated using an automatic calibration process in order to evaluate how far parameter values can
150 compensate for errors induced by grid coarsening. However, prior to the calibration, a sensitivity
151 analysis is performed for evaluating the influence of horizontal spatial discretization on parameter
152 sensitivities.

153 The sensitivity of each parameter included in the calibration process is evaluated using the composite
154 scaled sensitivity cSS_j (Hill, 1992; Anderman et al., 1996; Hill et al., 1998; Hill and Tiedeman, 2007):

$$155 \quad cSS_j = \left[\frac{\sum_{i=1}^{nobs} (dss_{ij})^2}{nobs} \Big|_b \right]^{1/2} \quad j = 1, npar \quad (1)$$

156 with the dimensionless scaled sensitivities dss_{ij} calculated as

$$157 \quad dss_{ij} = \frac{\partial y_i^{sim}}{\partial b_j} \Big|_b \times |b_j| \times w_{ii}^{1/2} \quad i = 1, nobs \quad j = 1, npar \quad (2)$$

158 The composite scaled sensitivity measures the information provided by the entire set of observations for
159 the estimation of the single parameter b_j . Large values correspond to sensitive parameters for which the
160 observations provide a lot of information. According to Hill et al. (1998) and Hill and Tiedeman (2007),
161 parameters with composite scaled sensitivities less than 1 or less than 0.01 of the largest composite
162 scaled sensitivity are poorly sensitive. Consequently, they could produce problems during the calibration
163 or calibrated parameters with large confidence intervals.

164 The calibration is performed using PEST (Doherty, 2005) enhanced with the temporary parameter
165 immobilization strategy developed by Skahill and Doherty (2006). The iterative local optimization
166 method implemented in PEST allows calculating the set of parameter values that produces the smallest
167 value of an objective function measuring the discrepancies between observed values and their
168 simulated equivalent. The objective function implemented in PEST is the weighted least-squares
169 objective function (L_2 norm):

$$170 \quad \Phi(\mathbf{b}) = \sum_{i=1}^{nobs} w_i \times [y_i^{obs} - y_i^{sim}(\mathbf{b})]^2 = \sum_{i=1}^{nobs} w_i \times r_i^2 \quad (3)$$

171 where $nobs$ is the number of observations of any kind, y_i^{obs} is the i^{th} observed value, $y_i^{sim}(\mathbf{b})$ is the
172 simulated equivalent to the i^{th} observed value calculated with the parameter values composing the
173 vector \mathbf{b} , w_i is the weight for the i^{th} contribution to the objective function, r_i is the i^{th} residual. However,
174 in presence of local minima in the objective function, this method based on local parameter sensitivities
175 does not always provide the set of parameter values corresponding to the global minimum. The use of
176 the temporary parameter immobilization strategy greatly reduces this eventuality. This strategy consists
177 in selectively withdrawing the most insensitive parameters from the estimation process when the
178 objective function improvement during a particular iteration is poor. This greatly heightens the capacity
179 of the estimation process to find the global minimum of the objective function. According to Doherty

180 (2005), calibration using truncated singular value decomposition, gives similar results since this method
181 also has the capacity of withdrawing insensitive parameters from the estimation process. Global
182 optimization methods ensuring to find the global minimum of the objective function are not used in this
183 study because they require a huge number of model runs which induces execution times tens or
184 hundreds of times longer than the execution times required by local optimization methods (Hill and
185 Tiedeman, 2007). This precludes using these methods for integrated surface and subsurface transient
186 flow simulations at large-scale due to their long execution times.

187 The set of parameters included in both the sensitivity analysis and the automatic calibration is
188 composed of 32 parameters corresponding to the physical state parameters found in the equations
189 representing surface and subsurface flow processes in HydroGeoSphere. The parameters present in the
190 equations representing the interception and evapotranspiration processes are not included. The set of
191 observations is composed of 24 discharge rates and 288 hydraulic head observations (1 per month and
192 per observation point for 2 years) produced with the synthetic catchment. The simulated values of
193 discharge and hydraulic head obtained with these calibrated models are saved for further graphical
194 model fit analysis and calculation of performance criteria.

195 **STEP 3 – Graphics of model fit, performance criteria, and linear confidence intervals on predictions.**

196 Graphical model fit analysis and calculation of performance criteria are performed for each model to
197 evaluate qualitatively and quantitatively the effects of spatial discretization on model performance and
198 to evaluate the improvement of model performance obtained with calibration. Additionally, the
199 influence of horizontal spatial discretization on model prediction uncertainties is evaluated using linear
200 confidence intervals on predictions.

201 **Graphics of model fit.** Graphical model fit analysis is somewhat subjective. However, it is good practice
202 to perform such a visual inspection prior to use numerical criteria for an objective evaluation of model

203 performance (Legates and McCabe, 1999; Hill and Tiedeman, 2007; Moriasi et al, 2007). Graphs
204 comparing observed and simulated values are the most widely used for evaluating model fit at a glance.
205 However, Hill and Tiedeman (2007) prefer using graphs such as weighted or unweighted simulated
206 values versus weighted residuals to facilitate the detection of model bias. If a model is unbiased, such
207 graphs exhibit weighted residuals evenly scattered about 0.0 for the entire range of values on the
208 horizontal axis. Weighted residuals wr_i are calculated as (Hill and Tiedeman, 2007):

$$209 \quad wr_i = w_i^{1/2} \times [y_i^{obs} - y_i^{sim}] = w_i^{1/2} \times r_i \quad (4)$$

210 The purpose of weighting is essentially to emphasize the most accurate observations. This is achieved by
211 specifying weights that are proportional or, preferably, equal to the inverse of the observation error
212 variances (Hill and Tiedeman, 2007):

$$213 \quad w_i = \frac{1}{\sigma_i^2} \quad (5)$$

214 where σ_i^2 is the true error variance of the i^{th} observation. Given these equations, in a graph of weighted
215 residuals versus unweighted simulated values, a cluster of negative weighted residuals indicate that
216 simulated values are systematically overestimated, and vice versa. Furthermore, with weights calculated
217 using a constant coefficient of variation, residuals are emphasized proportionally to their observed
218 value. Therefore, similar weighted residuals indicate similar relative errors. This way of emphasizing
219 residuals proportionally to their observed value is particularly useful for variables ranging over several
220 orders of magnitudes such as discharge.

221 **Performance criteria.** Performance criteria help quantifying model quality. They evaluate the level of
222 agreement between model and reality (Refsgaard and Henriksen, 2004). Typically, they depend on the
223 discrepancies between observed values and their simulated equivalent for a particular type of
224 observations (e.g. discharge or hydraulic heads). The performance criteria selected for this study are:

225 • The Nash-Sutcliffe efficiency criterion NSE_q (Nash and Sutcliffe, 1970):

$$226 \quad NSE_q = 1 - \frac{\sum_{t=1}^{nt} (q_t^{sim} - q_t^{obs})^2}{\sum_{t=1}^{nt} (q_t^{sim} - \mu^{obs})^2} \in] - \infty; 1] \quad (6)$$

227 where nt is the total number of timesteps, q_t^{sim} is the simulated discharge at timestep t , q_t^{obs} is the
 228 observed discharge at timestep t , and μ^{obs} is the mean of the observed values. If the simulated values
 229 perfectly match the observed values, $NSE_q = 1$. The lower the value of NSE_q , the poorer the model,
 230 negative values indicating that the mean observed value μ^{obs} gives a better description of the data than
 231 the simulated values q_t^{sim} . Weglarczyk (1998) and Gupta et al. (2009) suggest decomposing the Nash-
 232 Sutcliffe efficiency criterion for facilitating its interpretation. The decomposition of Gupta et. al (2009) is:

$$233 \quad NSE_q = 2 \times \frac{\sigma^{sim}}{\sigma^{obs}} \times r_{lin} - \left(\frac{\sigma^{sim}}{\sigma^{obs}} \right)^2 - \left(\frac{\mu^{sim} - \mu^{obs}}{\sigma^{obs}} \right)^2 \quad (7)$$

234 where r_{lin} is the linear correlation coefficient between q^{sim} and q^{obs} , σ^{sim} is the standard deviation of
 235 q^{sim} , μ^{sim} is the mean of q^{sim} , σ^{obs} is the standard deviation of q^{obs} , and μ^{obs} is the mean of q^{obs} . The
 236 first component uses the linear correlation coefficient for measuring the capacity of the model to
 237 reproduce the timing and the shape of the signal, the second component measures the capacity of the
 238 model to reproduce the standard deviations of the observations, and the third component measures the
 239 capacity of the model to reproduce the mean of the observations.

240 • The mass balance error MBE_q also known as bias, percent bias or relative bias (Gupta et al., 1999):

$$241 \quad MBE_q = \frac{\sum_{t=1}^{nt} (q_t^{sim} - q_t^{obs})}{\sum_{t=1}^{nt} q_t^{obs}} \times 100 = \frac{\mu^{sim} - \mu^{obs}}{\mu^{obs}} \times 100 \in] - 100; +\infty[\quad (8)$$

242 This performance criterion measures the tendency of the simulated values to be larger or smaller than
 243 their observed counterparts. If the fit is perfect, $MBE_q = 0$. If $MBE_q > 0$, simulated values are, on

244 average, greater than observed values, and vice versa. This performance criterion can also be used for
 245 hydraulic heads by substituting the observed and simulated discharges by the observed and simulated
 246 hydraulic heads in equation (7).

- 247 • The peak error PE_q (Aricò et al., 2009):

$$248 \quad PE_q = \left(\frac{q_{peak}^{sim}}{q_{peak}^{obs}} - 1 \right) \times 100 \in] - 100; +\infty[\quad (9)$$

249 where q_{peak}^{sim} is the simulated peak value, and q_{peak}^{obs} is the observed peak value. This performance
 250 criterion measures the capacity of the model to reproduce the peak in the hydrograph. If the observed
 251 peak is equal to the simulated peak, $PE_q = 0$. If $PE_q > 0$, the simulated peak is greater than the
 252 observed peak, and vice versa.

- 253 • The root mean squared error criterion RMS_h :

$$254 \quad RMS_h = \sqrt{\frac{1}{nt} \times \sum_{t=1}^{nt} (h_t^{sim} - h_t^{obs})^2} \in [0; +\infty[\quad (10)$$

255 where h_t^{sim} is the i^{th} simulated hydraulic head value, and h_t^{obs} is the i^{th} observed hydraulic head value.
 256 This performance criterion measures the discrepancies between observed hydraulic heads and their
 257 simulated equivalent for a particular observation point. If the simulated values perfectly match the
 258 observed values, $RMS_h = 0$. The greater the values, the poorer the model.

- 259 • The hydraulic head variations errors $HHVE_h$:

$$260 \quad HHVE_h = \left(\frac{h_{max}^{sim} - h_{min}^{sim}}{h_{max}^{obs} - h_{min}^{obs}} - 1 \right) \times 100 \in] - 100; +\infty[\quad (11)$$

261 where h_{max}^{sim} is the maximum simulated hydraulic head value, h_{min}^{sim} is the minimum simulated hydraulic
 262 head value, h_{max}^{obs} is the maximum observed hydraulic head value, and h_{min}^{obs} is the minimum observed
 263 hydraulic head value. This performance criterion is the counterpart of the peak error since it measures
 264 the capacity of the model to reproduce the magnitude of hydraulic head variations instead of measuring
 265 the capacity of the model to reproduce the peak in the hydrograph.

266 **Linear confidence intervals.** Linear and nonlinear confidence intervals help quantifying prediction
 267 uncertainties. Linear confidence intervals are calculated assuming that the model is linear in the vicinity
 268 of parameter values. They are not as accurate as nonlinear confidence intervals for nonlinear models.
 269 However, unlike nonlinear confidence intervals, linear confidence intervals only require trivial amount of
 270 execution time. Therefore, they are often the only confidence intervals calculable for physically-based
 271 and spatially-distributed models with lengthy execution times.

272 Linear confidence intervals on predictions have the form:

$$273 \quad z'_i \pm t_S(n, 1.0 - \frac{\alpha}{2}) \times s_{z'_i} \quad (12)$$

274 where z'_i is the i^{th} simulated prediction, $t_S(n, 1.0 - \frac{\alpha}{2})$ is the Student-t distribution with
 275 $n = (nobs - npar)$ and $\alpha = 0.05$ for 95% confidence intervals, and $s_{z'_i}$ is the standard deviation of the

276 prediction calculated as:

$$277 \quad s_{z'_i} = \left[\sum_{i=1}^{npar} \sum_{j=1}^{npar} \frac{\partial z'_i}{\partial b_j} \times V(\mathbf{b}) \times \frac{\partial z'_i}{\partial b_i} \right]^{1/2} \quad (13)$$

278 where $npar$ is the number of parameters, $\frac{\partial z'_i}{\partial b_j}$ is the sensitivity of the i^{th} prediction z'_i with respect to the
 279 j^{th} parameter b_j and $V(\mathbf{b})$ is the parameter variance-covariance matrix.

280 **3 Conceptual model**

281 The synthetic catchment is inspired by a real catchment located in the Condroz region of Belgium. This
282 region is characterized by a succession of limestone synclines and sandstone anticlines. The surface
283 materials of the synthetic catchment are assigned using a criterion combining elevation and slope
284 constraints. All surface materials are assigned a series of parameters required for simulating
285 interception, evapotranspiration, and surface flow processes. Appropriate values for these parameters
286 are extracted from the literature. They are listed in Appendix A. The subsurface materials of the
287 synthetic catchment are defined to represent the typical geology of the Condroz region: sandstones,
288 limestones, and shales constitute, from the crests to the center of the valley, the subsurface materials of
289 the synthetic catchment. Additionally, these formations are covered by alluvial deposits and loam. All
290 subsurface materials are assigned a series of parameters required for simulating subsurface flow
291 processes. Appropriate values for these parameters, including van Genuchten parameters governing
292 saturation-pressure relations in the vadose zone, are extracted from the literature. They are listed in
293 Appendix B. The synthetic catchment is illustrated in Figure 1.

294 The horizontal element size of the reference model progressively increases from 25 m near the surface
295 water network to 250 m far from the surface water network. The layer thickness progressively increases
296 from 1 m for the top layers corresponding to the vadose zone to 30 m for the bottom layers
297 corresponding to the saturated zone (5 layers of 1 m, 1 layer of 5 m, 1 layer of 10 m, and 1 layer of
298 30 m). The reference model is composed of 153,027 nodes and 269,872 elements. The grid of the
299 reference model is illustrated in Figure 2. Critical-depth boundary conditions are assigned to boundary
300 nodes of the surface domain. This type of boundary condition forces the water elevation at the
301 boundary to be equal to the water elevation for which the energy of the flowing water relatively to the
302 stream bottom is minimum (Therrien et al., 2012). No-flow boundary conditions are assigned to

303 boundary nodes of the subsurface domain. Water depths and hydraulic heads extracted from
304 preliminary simulations performed with the reference model are used as initial conditions for the
305 surface domain and the subsurface domain, respectively.

306 The set of observation points is constituted of 1 gauging station for discharge (G1) and 12 piezometers
307 evenly distributed in the synthetic catchment for hydraulic heads (Pz1 to Pz12). Two galleries (GAL1 and
308 GAL2) and four wells (W1 to W4) are used to simulate groundwater withdrawals. The set of observation
309 points and the galleries and wells are illustrated in Figure 1. As the models with a coarse grid are run
310 with monthly stress factors, discharge and hydraulic heads simulated at the observation points each day
311 of the 5-year reference simulation are monthly averaged for ensuring time consistency (Hill and
312 Tiedeman, 2007, p. 215). These monthly averaged discharge and hydraulic heads constitute the set of
313 reference observations used to calculate performance criteria for the simplified models. The reference
314 simulation is subdivided into warm-up, calibration, and validation periods. The warm-up is necessary for
315 obtaining simulated values independent of the initial conditions. Discharge and hydraulic heads
316 produced during the warm-up period are not included in the set of reference observations. Performance
317 criteria are only calculated for discharge and hydraulic heads produced during calibration and validation
318 periods. Linear confidence intervals on predictions are calculated for the validation period.

319 **4 Results and Discussion**

320 The models developed for evaluating the effects of spatial discretization on model performance and
321 model prediction uncertainties are referred as the *250 m*, *500 m*, *750 m*, and *1000 m* models. They are
322 characterized by a constant element size of 250 m, 500 m, 750 m, and 1000 m, respectively. As opposed
323 to the reference model, they are not refined near the surface water network. The purpose here consists
324 in evaluating the effects of ignoring such a refinement on the simulation of discharge and hydraulic
325 heads. Additionally, it also allows evaluating whether calibration can compensate for the errors induced

326 by ignoring such a refinement. As the reference model, they each have 8 layers (5 layers of 1 m, 1 layer
327 of 5 m, 1 layer of 10 m, and 1 layer of 30 m). The number of nodes, the number of elements, and the
328 execution times of the 250 m, 500 m, 750 m, and 1000 m models are presented in Table 1. The
329 comparison between the execution time of each model clearly shows the usefulness of coarsening grid
330 size for reducing the execution times.

331 **4.1 Comparison of model performance before calibration**

332 Graphs of model fit and performance criteria are used together for comparing the performance of the
333 250 m, 500 m, 750 m, and 1000 m models run with the same parameter values than those used in the
334 reference model i.e. without any calibration.

335 Graphs comparing reference values of discharge and hydraulic heads produced with the reference
336 model and their simulated equivalent obtained with the models with a coarse grid indicate that
337 discharge is most often underestimated during low flow periods and overestimated during high flow
338 periods (Figure 3-A). The underestimation is almost identical for each model. The overestimation is
339 higher for models with a coarse horizontal spatial discretization. This is clearly visible on peak discharge.

340 Graphs of unweighted simulated values versus weighted residuals support these findings. These graphs
341 particularly highlight the underestimation of discharge by each model during low flow periods and the
342 overestimation of discharge during high flow rates by the coarsest ones (Figure 3-B).

343 The influence of horizontal spatial discretization on hydraulic head simulation is less visible (Figure 3-A).
344 However, weighted residuals are in general greater for the coarsest models (Figure 3-B). This shows that
345 the simulation of hydraulic heads is poorer with the coarsest models.

346 Graphical model fit analysis is confirmed by performance criteria. As the grid is coarsened, NSE_q values
347 tend to decrease and RMS_h values tend to increase (Figure 4-A). This indicates that simulation of both

348 discharge and hydraulic heads is deteriorated. For discharge, Gupta's decomposition of NSE_q shows
349 that the standard deviation of discharge is overestimated by the coarsest models (Table 2). This is visible
350 to the greater values of Gupta's second terms. This is also supported by the increasing values of PE_q^{yr1}
351 and PE_q^{yr2} showing that peak discharge, and so the standard deviation of the hydrograph, are
352 overestimated by the coarsest models (Table 2). Gupta's decomposition also shows that the 250 m
353 model lacks to properly simulate the average magnitude of discharge. This is why NSE_q value for this
354 model is lower than NSE_q value for the 500 m model. This is confirmed by the values of MBE_q which
355 shows that the 250 m model underestimates the average magnitude of discharge by almost 15%. This is
356 related to the fact that the underestimation of discharge during low flow periods is not compensated by
357 the overestimation of discharge during high flow rates as it is the case for the other models. For
358 hydraulic heads, the absolute values of MBE_h are in general low for each model (Table 3). This indicates
359 that models are not significantly biased in terms of hydraulic heads. However, the range of MBE_h values
360 is in general wider for the coarsest models. Although the absolute values of $HHVE_h^{yr1}$ and $HHVE_h^{yr2}$
361 are in general greater for the coarsest models, the ranges of $HHVE_h^{yr1}$ and $HHVE_h^{yr2}$ are similar for
362 each model (Table 3).

363 The comparison of model performance performed in this section indicates that coarsening the grid
364 mainly deteriorates the simulation of discharge. Common to each model tested, the underestimation of
365 discharge during low flow periods is due to a poor representation of the surface water network which
366 precludes from properly simulating groundwater-surface water interactions that constitute the key
367 component of the hydrograph during dry seasons. As previously mentioned, this problem of poor
368 representation of the surface water network is also mentioned by Refsgaard (1997) and Vázquez et al.
369 (2002). The overestimation of discharge by the coarsest models during high flow periods is related to

370 the use of large elements which induces a smoothing of surface slopes and facilitates runoff, especially
371 during wet seasons. The object of the next section is to evaluate how calibration can compensate for the
372 errors induced by coarsening the grid.

373 **4.2 Comparison of model performance after calibration**

374 A sensitivity analysis is performed for each parameter prior to the calibration. The composite scaled
375 sensitivities calculated on the calibration period (24 discharge observations and 288 hydraulic head
376 observations) for each parameter included in the calibration (32 parameters) are illustrated in Figure 5.
377 Whatever the spatial discretization, the ranking of the most sensitive parameters and the magnitude of
378 the composite scaled sensitivities are almost identical. This suggests that parameter sensitivities are not
379 highly dependent on the grid size. The most sensitive parameter is always the van Genuchten
380 parameters β_{VG} of Mat I – loam and Mat II – alluvial deposits (top layers of the models). This parameter,
381 related to the pore-size distribution in the porous medium, defines the shape of the water retention
382 curve. The other most sensitive parameters are the hydraulic conductivity K of Mat IV – limestones 2,
383 probably because most of the observation points are located in this material, and the van Genuchten
384 parameter β_{VG} of Mat IV – limestones 1, Mat V – limestones 2 and Mat VI – sandstones. The van
385 Genuchten parameter α_{VG} of Mat IV – limestones 1, Mat V – limestones 2 and Mat VI – sandstones as
386 well as the hydraulic conductivity K of Mat I – loam have also a relatively high sensitivity. The fact that
387 van Genuchten parameters, especially the parameter β_{VG} of the materials constituting the top layers of
388 the models, are systematically among the most sensitive parameters suggests that fully-integrated and
389 physically-based models are highly sensitive to parameters governing the infiltration process in the
390 vadose zone and the groundwater recharge.

391 The improvement of model performance with calibration with PEST is evaluated using the same graphs
392 of model fit and the same performance criteria than in the previous section. Graphs of model fit show

393 that calibration significantly improves the simulation of discharge and, to a lesser extent, hydraulic
394 heads for each model (Figure 6-A). Additionally, after calibration, weighted residuals are almost
395 randomly distributed which suggests that calibrated models are less biased (Figure 6-B). Performance
396 criteria support these findings since NSE_q and RMS_h values are significantly greater and lower,
397 respectively, after calibration (Figure 4-B). The values of Gupta's terms together with the values of
398 MBE_q , $PE_q^{yr 1}$, and $PE_q^{yr 2}$ calculated for the calibrated models indicate that both the mean and the
399 standard deviation of flow rates are better simulated (Table 4). The improvement of hydraulic head
400 simulation is not so clear. When observed and simulated hydraulic heads are shifted, the calibration
401 process strives for reducing this systematic error. Therefore, the improvement of average hydraulic
402 head magnitudes is sometimes obtained to the detriment of the improvement of hydraulic head
403 variations. This is why the absolute values and the range of MBE_h are most often lower than those
404 obtained with the models before calibration, while the absolute values and the ranges of $HHVE_h^{yr 1}$ and
405 $HHVE_h^{yr 2}$ are identical or even greater than those obtained with the models before calibration (Table
406 5). This shows that calibration has limitations. Furthermore, although most of them are still within
407 reasonable ranges, some calibrated parameter values are far from their values in the reference model
408 (Table 6). Such an observation is only possible for synthetic catchments for which reference parameter
409 values are exactly known. The only verification possible for real catchments consists in making sure that
410 calibrated parameter values are plausible with regards to field or laboratory data. However, as shown by
411 Brunner et al. (2012b), accurately evaluating certain combinations of parameters can be sufficient to
412 produce predictions with a good level of confidence, which means that it is not always necessary to
413 accurately evaluate each parameter individually. Therefore, in spite of its limitations, calibration is
414 essential for improving model performance, either inside or outside the calibration period. As illustrated
415 in Figure 7, calibration indeed leads to greater values of NSE_q and lower values of RMS_h also during the

416 validation period. The object of the next section is to evaluate whether grid coarsening leads to greater
417 model prediction uncertainties.

418 **4.3 Comparison of model prediction uncertainties**

419 Linear confidence intervals on predictions are calculated for discharge and hydraulic heads simulated in
420 the validation period with the calibrated models. They are illustrated in Figures 8 to 11. The linear
421 confidence intervals calculated for discharge are almost identical for the *250 m*, *500 m*, and *750 m*
422 models. They are even sometimes narrower for the *500 m* or the *750 m* models than for the *250 m*
423 model. However, especially for high flow periods, they are far wider for the *1000 m* model. The linear
424 confidence intervals calculated for hydraulic heads are quite similar for each model and once more the
425 narrowest intervals are not always obtained for the *250 m* model.

426 The analysis of model prediction uncertainties indicates that coarsening model grid mainly influences
427 the uncertainties on discharge predictions. This is not surprising since the comparison of model
428 performance shows that discharge simulation is more sensitive to grid size than hydraulic head
429 simulation. However, the uncertainties on discharge predictions significantly increase only for a very
430 coarse grid and even if graphs of model fit and performance criteria suggest that the model is good.
431 Therefore, to some extent, it is possible to simplify a model by coarsening its grid without increasing
432 model prediction uncertainties. This is consistent with the study of Brunner et al. (2012b) focusing on
433 parameter identifiability and predictive uncertainty. This study highlights the sliding nature of
434 complexity versus simplicity and shows that predictive power may lose little if the model is simplified
435 appropriately.

436 **4.4 Guidelines for selection of a proper horizontal spatial discretization for large-scale flow models**

437 synthetic catchment can always be considered as far from reality. Therefore, caution should be
438 exercised when using results of this study for selecting a proper horizontal spatial discretization for a

439 given site-specific study. However, a series of general guidelines can be drawn from this study. As an
440 example, in the framework of use of paired simple and complex models to reduce predictive bias and
441 quantify uncertainty (Doherty and Christensen, 2011), these guidelines could be used for helping
442 modelers selecting a proper horizontal spatial discretization for the simple model.

443 Large-scale physically-based and spatially-distributed model development consists in finding a
444 compromise between model accuracy and model portability i.e. maximizing model performance and
445 minimizing prediction uncertainty while limiting the execution times. Given the results of this study, for
446 catchments of a few hundreds square kilometer, an element size of the order of 500 m is the best
447 compromise for obtaining good model performance with tractable execution times without significantly
448 increasing prediction uncertainty. With a coarser horizontal spatial discretization, the relative reduction
449 of execution times is limited with respect to the probability of increasing prediction uncertainty. With a
450 finer horizontal spatial discretization, the execution times strongly increase without any significant
451 reduction of prediction uncertainty.

452 **5 Summary and Conclusions**

453 The present study focuses on the effects of horizontal spatial discretization on large-scale flow model
454 performance and model prediction uncertainties using a fully-integrated hydrological model of a
455 synthetic catchment. This kind of large-scale fully-integrated hydrological model is increasingly used in
456 water management for predicting the evolution of both the integrated response (discharge) and the
457 distributed response (hydraulic heads) of catchments. However, these models are characterized by
458 lengthy execution times and model grids are often coarsened for reducing these execution times.
459 Therefore, it is crucial to evaluate the influence of such a grid coarsening on model performance and
460 model prediction uncertainties. This study shows that:

- 461 • Grid coarsening mainly influences the simulation of discharge with an underestimation of
462 discharge during low flow periods and a progressive overestimation of peak discharge as
463 horizontal spatial discretization is coarsened. This is related to a poor representation of the
464 surface water network and the smoothing of surface slopes that prevent from properly
465 simulating surface water-groundwater interactions and runoff process.
- 466 • Parameter sensitivities are not significantly influenced by grid coarsening and model errors
467 induced by grid coarsening can be compensated by calibration (preferably using an automatic
468 calibration process). Furthermore, calibration improves model performance either inside or
469 outside the calibration period. However, calibration has limitations and model errors are
470 potentially compensated at the cost of less plausible parameter values.
- 471 • Grid coarsening mainly influences the uncertainty on discharge predictions. However, model
472 prediction uncertainties on discharge only increase significantly for very coarse horizontal
473 spatial discretizations.

474 As uncertainty analyses have become essential in natural system modeling, this is encouraging since grid
475 coarsening greatly reduces execution times and such analyses can only be performed for model with
476 relatively short execution times.

477 **References**

478 Anderman, E., Hill, M.C., Poeter, E.P., 1996. Two-dimensional advective transport in ground-water flow
479 parameter estimation. *Ground Water* 34(6), 1001-1009.

480 Aricò, C., Nasello, C., Tucciarelli, T., 2009. Using steady-state water level data to estimate channel
481 roughness and discharge hydrograph. *Advances in Water Resources* 32(8), 1223-1240.

482 Bauer, S., Beyer, C., Kolditz, O., 2006. Assessing measurement uncertainty of first-order degradation
483 rates in heterogeneous aquifers. *Water Resources Research* 42(W01420).

484 Beyer, C., Bauer, S., Kolditz, O., 2006. Uncertainty assessment of contaminant plume length estimates in
485 heterogeneous aquifers. *Journal of Contaminant Hydrology* 87(1-2), 73-95.

486 Brunner, P., Doherty, J., Simmons, C.T., 2012a. Uncertainty assessment and implications for data
487 acquisition in support of integrated hydrologic models. *Water Resources Research* 48(W07513).

488 Brunner, P., Simmons, C.T., 2012b. HydroGeoSphere : A fully-integrated, physically-based hydrological
489 model. *Ground Water* 50(2), 170-176.

490 Doherty, J., 2005. PEST – Model-independent parameter estimation – User Manual – 5th edition.
491 Watermark Numerical Computing.

492 Doherty, J., Christensen, S., 2011. Use of paired simple and complex models to reduce predictive bias and
493 quantify uncertainty. *Water Resources Research* 47(W12534).

494 Downer, C.W., Ogden, F.L., 2004. Appropriate vertical discretization of Richards' equation for two-
495 dimensional watershed-scale modelling. *Hydrological Processes* 18, 1-22.

496 Ebel, B., Loague, K., 2006. Physics-based hydrologic-response simulation: seeing through the fog of
497 equifinality. *Hydrological Processes* 20(13), 2887-2900.

498 Frei, S., Fleckenstein, J.H., Kollet, S.J., Maxwell, R.M., 2009. Patterns and dynamics of river-aquifer
499 exchange with variably-saturated flow using a fully-coupled model. *Journal of Hydrology* 375, 383-393.

500 Goderniaux, P., Brouyère, S., Fowler, H., Blenkinsop, S., Therrien, R., Orban Ph., Dassargues, A., 2009.
501 Large scale surface-subsurface hydrological model to assess climate change impacts on groundwater
502 reserves. *Journal of Hydrology* 373(1-2), 122-138.

503 Goderniaux, P., Brouyère, S., Blenkinsop, S., Burton, A., Fowler, H., Orban Ph., Dassargues, A., 2011.
504 Modeling climate change impacts on groundwater resources using transient stochastic climatic
505 scenarios. *Water Resources Research* 47 (W12516).

506 Gupta, H., Sorooshian, S., Yapo, P., 1999. Status of automatic calibration for hydrologic models:
507 comparison with multilevel expert calibration. *Journal of Hydrologic Engineering* 4(2), 135-143.

508 Gupta, H., Kling, H., Yilmaz, K., Martinez, G., 2009. Decomposition of the mean squared error and NSE
509 performance criteria: implications for improving hydrological modelling. *Journal of Hydrology*, 377(1-2),
510 80-91.

511 Hill, M.C., 1992. A computer program (MODFLOWP) for estimating parameters of a transient, three-
512 dimensional, ground-water flow model using nonlinear regression. Open-File Report 91-484, U.S.
513 Geological Survey.

514 Hill, M.C., Cooley, R., Pollock, D., 1998. A controlled experiment in groundwater flow model calibration
515 using nonlinear regression. *Ground Water* 36(3), 520-535.

516 Hill, M.C., Tiedeman, C.R., 2007. Effective groundwater model calibration with analysis of sensitivities,
517 predictions and uncertainty. Hoboken, NJ, USA: John Wiley & Sons.

518 Irvine, D.J., Brunner, P., Hendricks Franssen, H.-J., Simmons C.T., 2012. Heterogeneous or
519 homogeneous? Implications of simplifying heterogeneous streambeds in models of losing streams.
520 *Journal of Hydrology* 424-425, 16-23.

521 Kristensen, K.J., Jensen, S.E., 1975. A model for estimating actual evapotranspiration from potential
522 evapotranspiration. *Nordic Hydrology* 6, 170-188.

523 Legates, D., McCabe, G., 1999. Evaluating the “goodness-of-fit” measures in hydrologic and
524 hydroclimate model validation. *Water Resources Research* 35(1), 233-241.

525 Li, Q., Unger, A., Sudicky, E., Kassenaar, D., Wexler, E., Shikaze, S., 2008. Simulating the multi-seasonal
526 response of a large-scale watershed with a 3D physically-based hydrologic model. *Journal of Hydrology*
527 357(3-4), 317-336.

528 Meyerhoff, S.B., Maxwell, R.M., 2011. Quantifying the effects of subsurface heterogeneity on hillslope
529 runoff using a stochastic approach. *Hydrogeology Journal* 19, 1515-1530.

530 Moriasi, D., Arnold, J., Van Liew, M., Binger, R., Harmel, R., Veith, T., 2007. Model evaluation guidelines
531 for systematic quantification of accuracy in watershed simulations. *Transactions of the ASABE* 50(3),
532 885-900.

533 Nash, J., Sutcliffe, J., 1970. River flow forecasting through conceptual models I – A discussion of
534 principles. *Journal of Hydrology* 10(3), 282-290.

535 Poeter, E.P., McKenna, S., 1995. Reducing uncertainty associated with ground-water flow and transport
536 predictions. *Ground Water* 33(6), 899-904.

537 Poeter, E.P., Hill, M.C., 1997. Inverse models: a necessary next step in ground-water modeling. *Ground*
538 *Water* 35(2), 250-260.

539 Refsgaard, J.C, 1997. Parameterisation, calibration and validation of distributed hydrological models.
540 *Journal of Hydrology* 198, 69-97.

541 Refsgaard, J.C., Henriksen, H., 2004. Modelling guidelines – terminology and guiding principles.
542 *Advances in Water Resources* 27(1), 71-82.

543 Schäfer, D., Schlenz, B., Dahmke, A., 2004. Evaluation of exploration and monitoring methods for
544 verification of natural attenuation using the virtual aquifer approach. *Biodegradation* 15(6), 453-465.

545 Sciuto, G., Diekkrüger, B., 2010. Influence of soil heterogeneity and spatial discretization on catchment
546 water balance modeling. *Vadose Zone Journal* 9, 955-969.

547 Skahill, B., Doherty, J., 2006. Efficient accommodation of local minima in watershed model calibration.
548 *Journal of Hydrology* 329(1-2), 122-139.

549 Therrien, R., McLaren, R., Sudicky, E., Park, Y.-J., 2012. *HydroGeoSphere – A three-dimensional*
550 *numerical model describing fully-integrated subsurface and surface flow and solute transport. Manual,*
551 *Groudwater Simulations Group.*

552 Vázquez R.F., Feyen, L., Feyen, J., Refsgaard, J.C., 2002. Effect of grid size on effective parameters and
553 model performance of the MIKE-SHE code. *Hydrological Processes* 16, 355-372.

554 Weglarczyk, S., 1998. The interdependence and applicability of some statistical quality measures for
555 hydrological models. *Journal of Hydrology* 206(1-2), 98-103.

556

557

558 **Appendices**

559 **Appendix A: Interception, evapotranspiration and surface flow parameters used in the synthetic case**

Interception and evapotranspiration	Mat 1	Mat 2	Mat 3	Mat 4
LAI [-]	0.40	3.53	5.12	-
L_r [m]	0	2.30	2.90	-
L_e [m]	2.00	2.00	2.00	-
θ_{e2} & θ_{t2} [-]	0.60	0.60	0.60	-
θ_{e1} & θ_{t1} [-]	0.96	0.96	0.96	-
C₁ [-]	0.31	0.31	0.31	-
C₂ [-]	0.15	0.15	0.15	-
C₃ [-]	10.00	10.00	10.00	-
C_{int} [m]	5.00 × 10 ⁻⁵	5.00 × 10 ⁻⁵	5.00 × 10 ⁻⁵	-
S⁰_{int} [m]	0	0	0	-
Surface flow	Mat 1	Mat 2	Mat 3	Mat 4
n_{xx} & n_{yy} [m^{-1/3}s]	0.012	0.200	0.600	0.025
H_{sto} [m]	0.002	0.002	0.002	0.002
L_c [m]	1.00 × 10 ⁻¹	1.00 × 10 ⁻¹	1.00 × 10 ⁻¹	1.00 × 10 ⁻¹

560 LAI = Leaf Area Index; L_r = root depth; L_c = evaporation depth; θ_{e2} & θ_{t2} and θ_{e1} & θ_{t1} = evaporation and
 561 transpiration limiting saturations; C₁, C₂, and C₃ = transpiration fitting parameters; C_{int} = canopy storage
 562 parameter; S⁰_{int} = initial interception storage; n_{xx} & n_{yy} = Manning roughness coefficients; H_{sto} = rill
 563 storage height; L_c = coupling length.

564 Parameter values are extracted from the literature:

- 565 • for parameters related to interception and evapotranspiration processes, see Andersen et al.,
566 2002; Asner et al., 2003; Canadell et al., 1996; Dickinson et al., 1991; Goderniaux, 2010; Graham
567 and Kilde, 2002; Islam, 2004; Kristensen and Jensen, 1975; Li et al., 2008; Panday and Huyakorn,
568 2004; Schroeder et al., 2004; Therrien et al., 2005; Vázquez et al., 2002; Vázquez and Feyen,
569 2003.
- 570 • for parameters related to surface flow processes, see Brutsaert, 2005; Fetter, 2001; Hornberger
571 et al., 1998; Jones, 2005; Li et al., 2008; McCuen, 1989.

572 **Appendix B: Subsurface flow parameters used in the synthetic case**

Subsurface flow	Mat I	Mat II	Mat III	Mat IV	Mat V	Mat VI
K [ms⁻¹]	5.00×10^{-7}	1.00×10^{-6}	1.00×10^{-5}	1.00×10^{-4}	2.50×10^{-4}	5.00×10^{-5}
S_s [m⁻¹]	1.00×10^{-4}	1.00×10^{-4}	1.00×10^{-4}	1.00×10^{-4}	1.00×10^{-4}	1.00×10^{-4}
θ_s [-]	4.10×10^{-1}	4.10×10^{-1}	2.50×10^{-2}	1.00×10^{-1}	1.00×10^{-1}	7.50×10^{-2}
S_{wr} [-]	9.76×10^{-2}	9.76×10^{-2}	0	0	0	0
α_{vG} [m⁻¹]	2.67	2.67	6.08×10^{-12}	3.65×10^{-2}	3.65×10^{-2}	3.65×10^{-2}
β_{vG} [-]	1.45	1.45	0.62	1.83	1.83	1.83
γ_{vG} [-]	1-1/β _{vG}	1-1/β _{vG}	38,671.00	1-1/β _{vG}	1-1/β _{vG}	1-1/β _{vG}

573 K = saturated hydraulic conductivity; S_s = specific storage; θ_s = saturated water content; S_{wr} = residual
574 water saturation; α_{vG}, β_{vG}, and γ_{vG} = van Genuchten parameters.

575 Parameter values are extracted from the literature:

- 576 • for parameters related to subsurface flow processes, see Brouyère et al., 2009; Freeze and
577 Cherry, 1979; Jones, 2005; Radcliffe, 2000; Ramos da Silva et al., 2008; Roulier et al., 2006.

578

579

580 **Figure captions**

581 **Figure 1** The reference model is assigned surface materials depending on elevation and slope constraints
582 and subsurface materials following the typical syncline structure of catchments located in the Condroz
583 region of Belgium. A gauging station (G1) and twelve piezometers (Pz1 to Pz12) are used to obtain
584 reference observations in terms of discharge and hydraulic heads, respectively. Two galleries (GAL1 and
585 GAL2) and four wells (W1 to W4) are used to simulate groundwater withdrawals.

586 **Figure 2** The grid of the reference model is refined horizontally (element side length from 25 m to
587 250 m) and vertically (layer thickness from 1 m to 30 m). The total number of nodes is 153,027.

588 **Figure 3 A.** As spatial discretization gets coarser, discharge simulation and, to a lesser extent, hydraulic
589 head simulation is progressively deteriorated. **B.** While each model underestimates discharge during low
590 flow periods, discharge during high flow periods is only overestimated by the coarsest models. This is
591 highlighted by the graphs of weighted residuals.

592 **Figure 4 A.** As horizontal spatial discretisation gets coarser, NSE_q values are in general lower and RMS_h
593 values are in general higher, this indicates that the simulation of both discharge and hydraulic heads are
594 progressively deteriorated. **B.** The higher values of NSE_q and the lower values of RMS_h obtained with the
595 calibrated models indicate that calibration significantly improves the simulation of both discharge and
596 hydraulic heads.

597 **Figure 5** Whatever the spatial discretization, the ranking of the most sensitive parameters and the
598 magnitude of the composite scaled sensitivities are almost similar. This suggests that parameter
599 sensitivities are not highly dependent on the grid size.

600 **Figure 6** Calibration significantly improves model performance even for the coarsest models. **B.**
601 Weighted residuals obtained with the calibrated models are almost randomly distributed. This indicates
602 that calibration reduces model bias.

603 **Figure 7** Values of NSE_q and RMS_h calculated for the validation period indicate that calibration also
604 improves model performance outside the calibration period.

605 **Figure 8** The 95% linear confidence intervals calculated for discharge only increase significantly for the
606 coarsest model.

607 **Figures 9 to 11** The 95% linear confidence intervals calculated for hydraulic heads are quite similar for
608 each model.

609

610

611 **Table captions**

612 **Table 1** Comparison of the number of nodes, number of elements, and execution times of the 250 m,
613 500 m, 750 m, and 1000 m models. The gain in execution time is tremendous when element size is
614 increased.

615 **Table 2** Values of NSE_q , MBE_q , $PE_q^{yr 1}$, and $PE_q^{yr 2}$ calculated for the 250 m, 500 m, 750 m, and 1000 m
616 models. When spatial discretisation gets coarser, the variance of the hydrograph is poorly simulated
617 (Gupta's 2nd term).

618 **Table 3** Values of MBE_h , $HHVE_h^{yr 1}$, and $HHVE_h^{yr 2}$ calculated for the 250 m, 500 m, 750 m, and 1000 m
619 models.

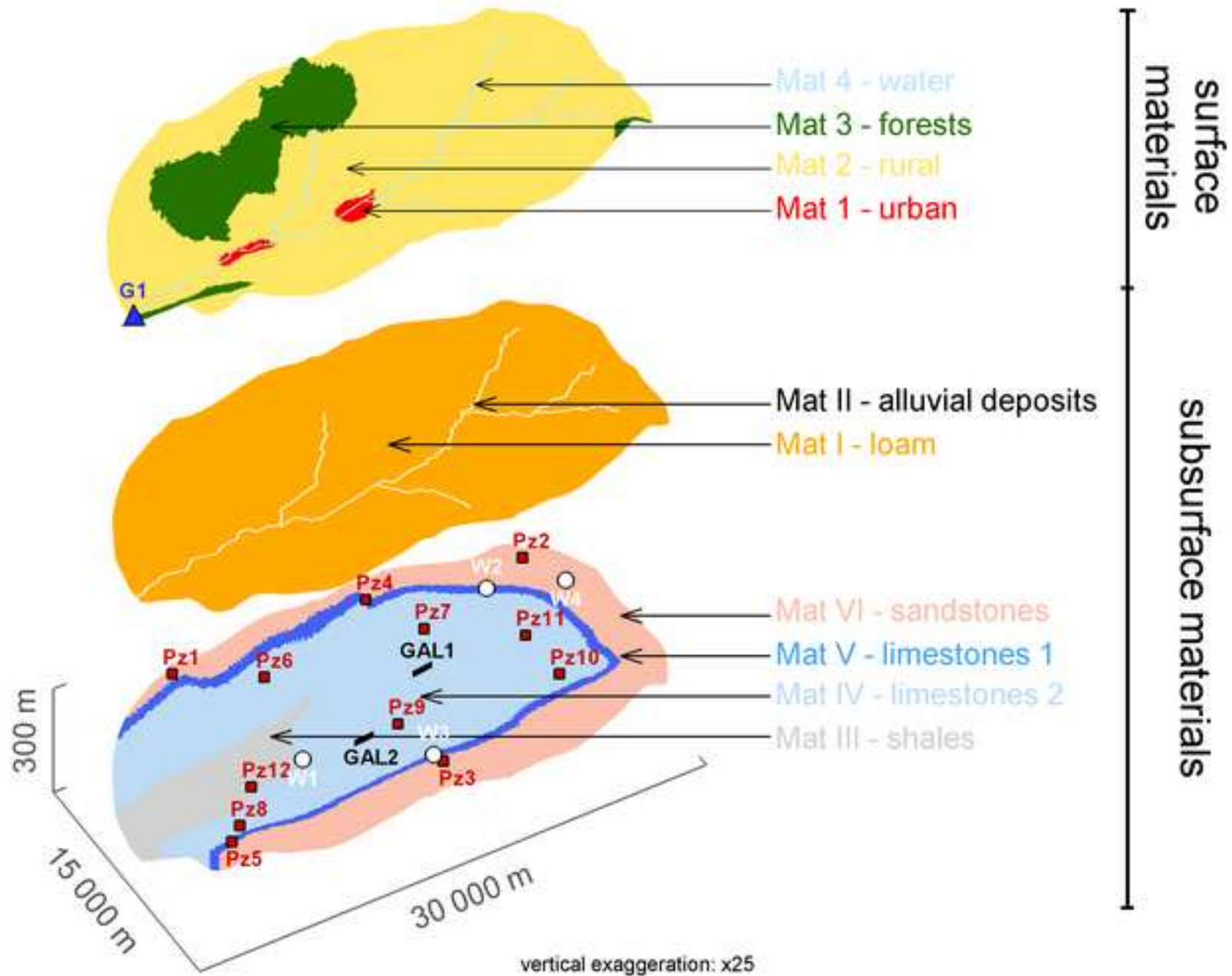
620 **Table 4** Values of NSE_q , MBE_q , $PE_q^{yr 1}$, and $PE_q^{yr 2}$ calculated for the calibrated 250 m, 500 m, 750 m,
621 and 1000 m models. Values in green are improved with regards to the corresponding models before
622 calibration. Values in red are deteriorated with regards to the corresponding forward models.

623 **Table 5** Values of MBE_h , $HHVE_h^{yr 1}$, and $HHVE_h^{yr 2}$ calculated for the calibrated 250 m, 500 m, 750 m,
624 and 1000 m models. Values in green are improved with regards to the corresponding models without
625 calibration. Values in red are deteriorated with regards to the corresponding forward models.

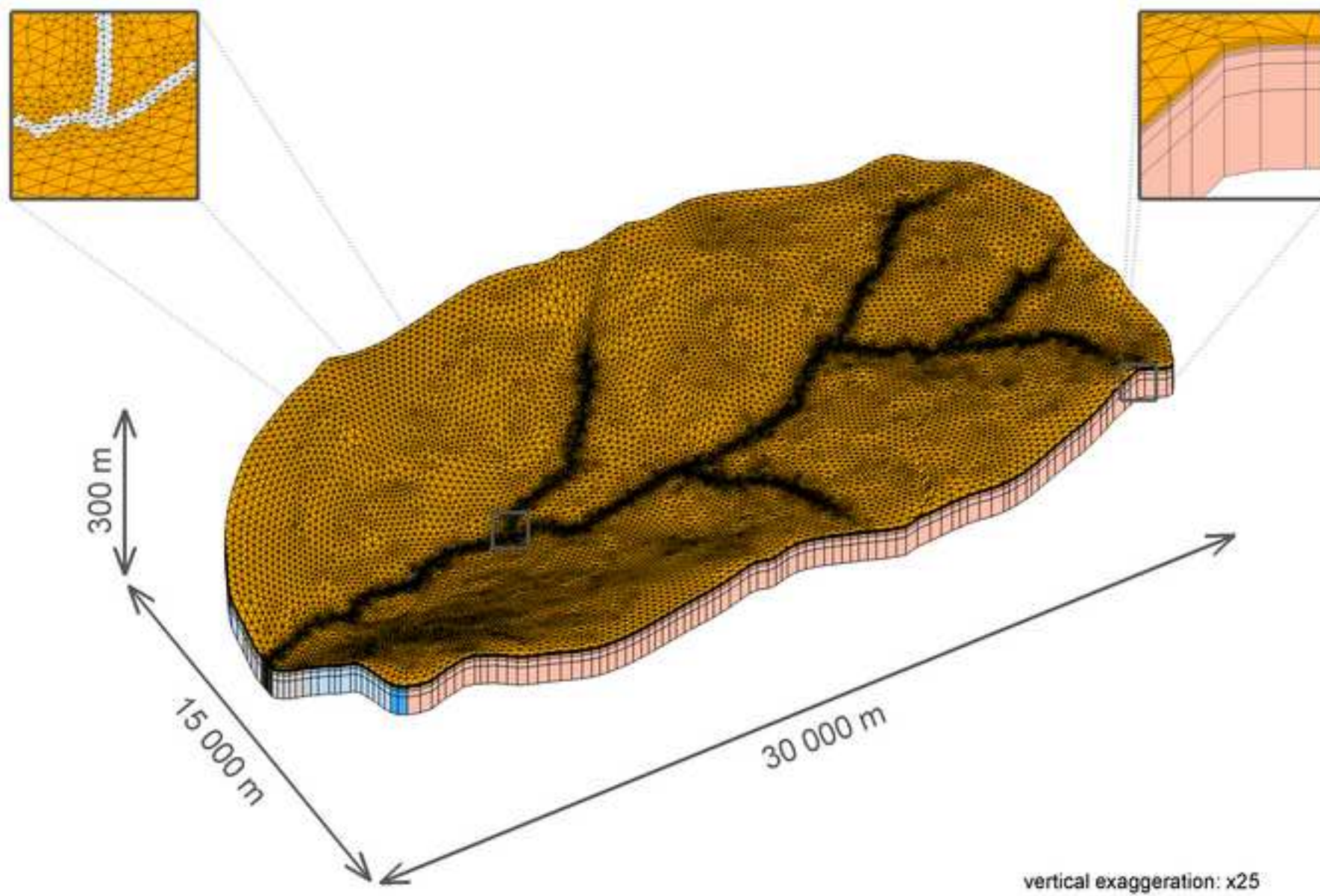
626 **Table 6** Comparison of the reference value of the most sensitive parameters and their value after
627 calibration.

Figure_1

[Click here to download high resolution image](#)



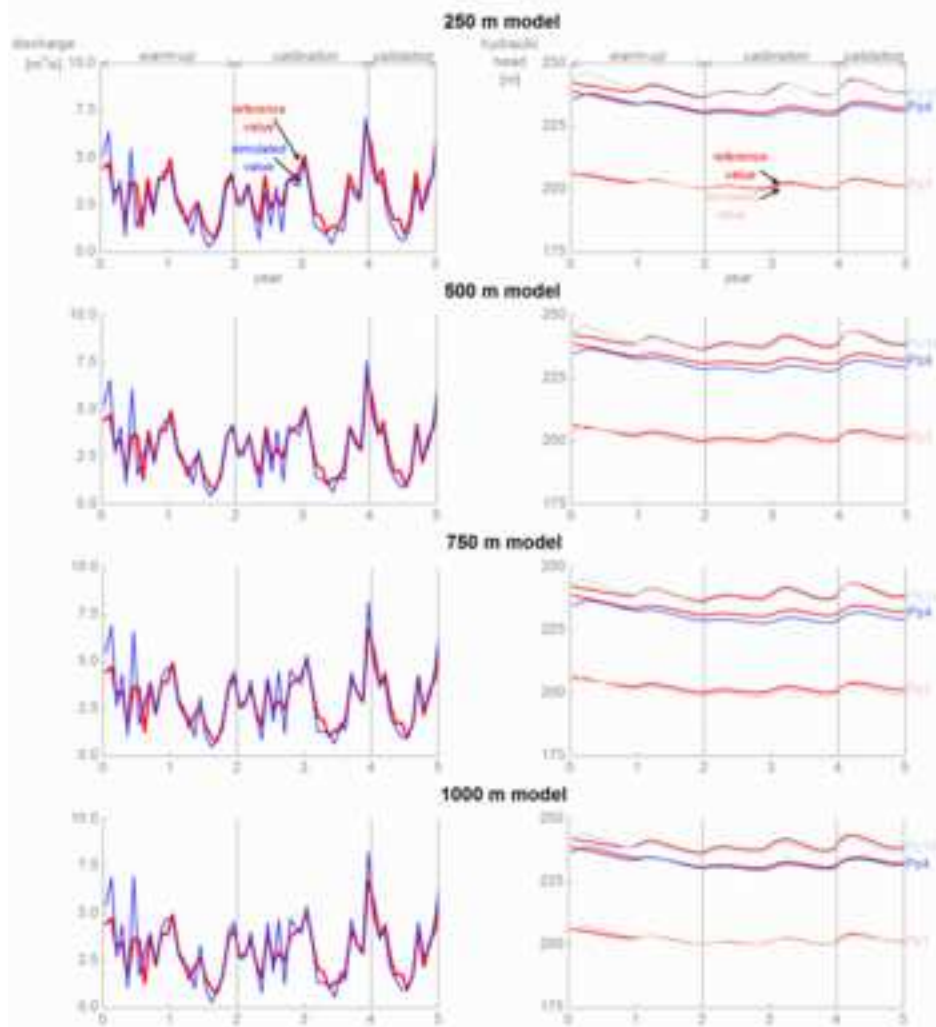
Figure_2
[Click here to download high resolution image](#)



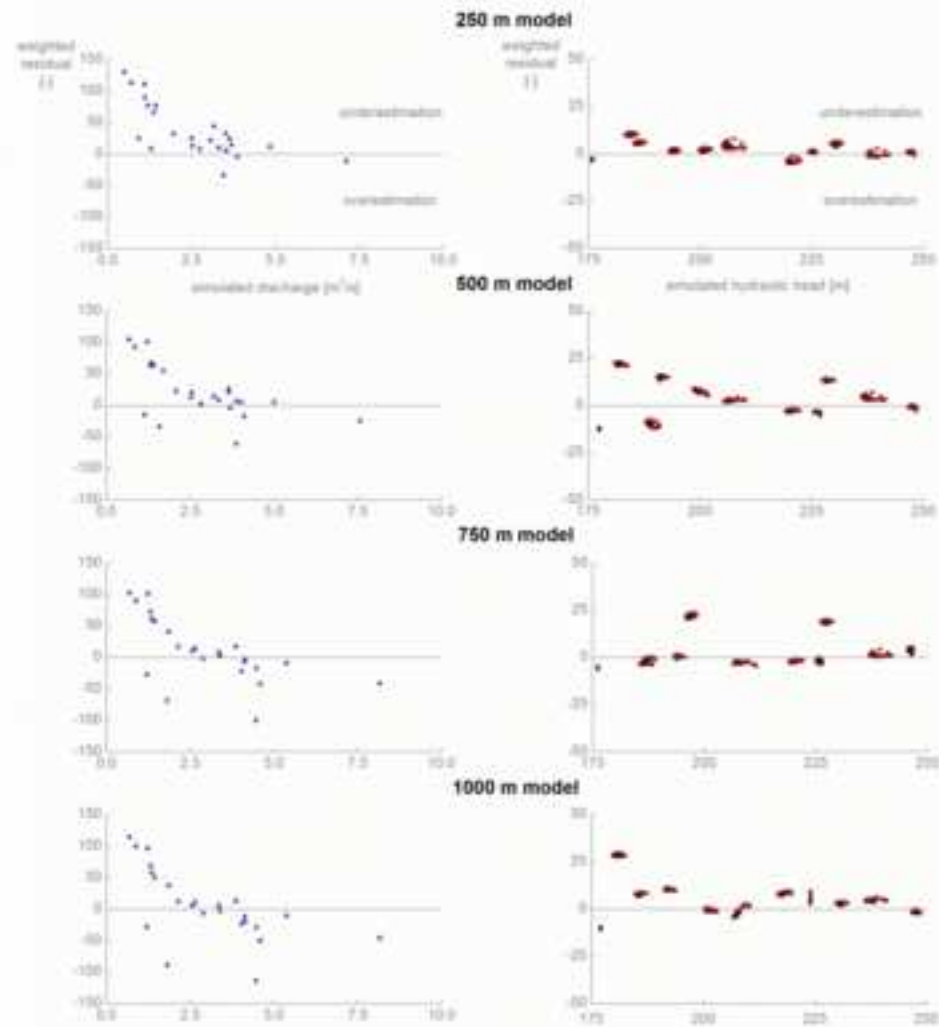
Figure_3

[Click here to download high resolution image](#)

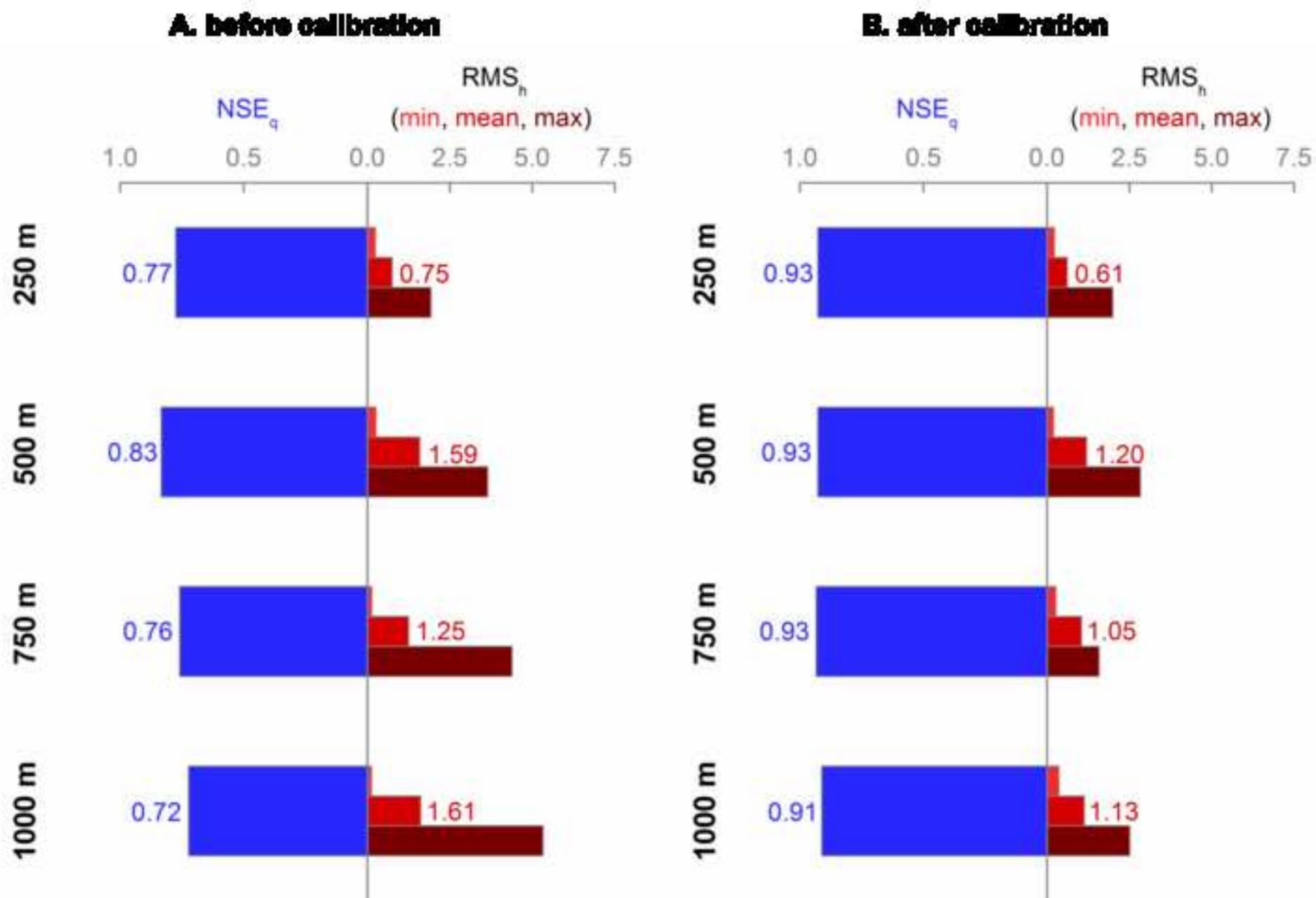
A. reference vs simulated values



B. simulated values vs weighted residuals

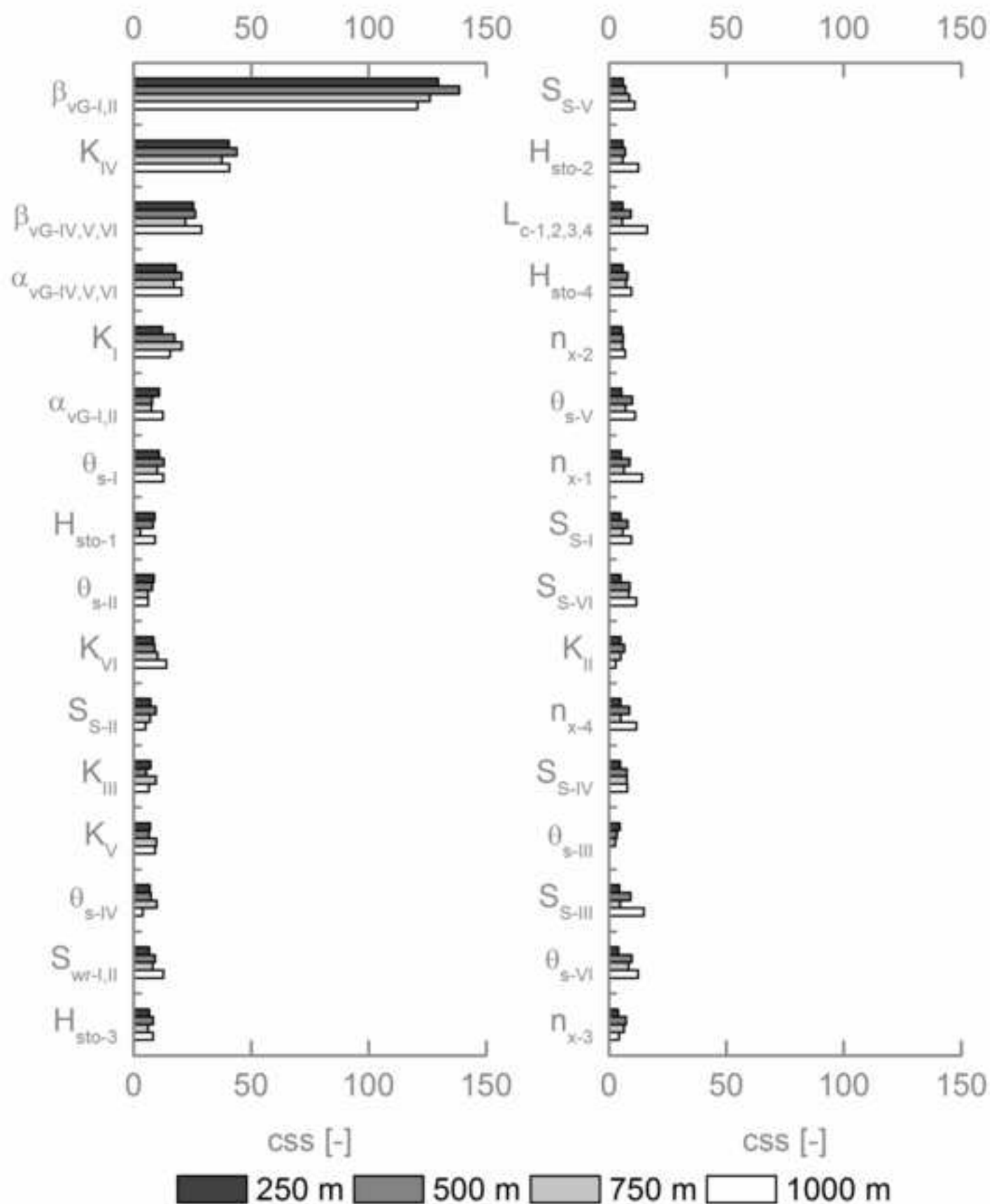


Figure_4

[Click here to download high resolution image](#)

Figure_5

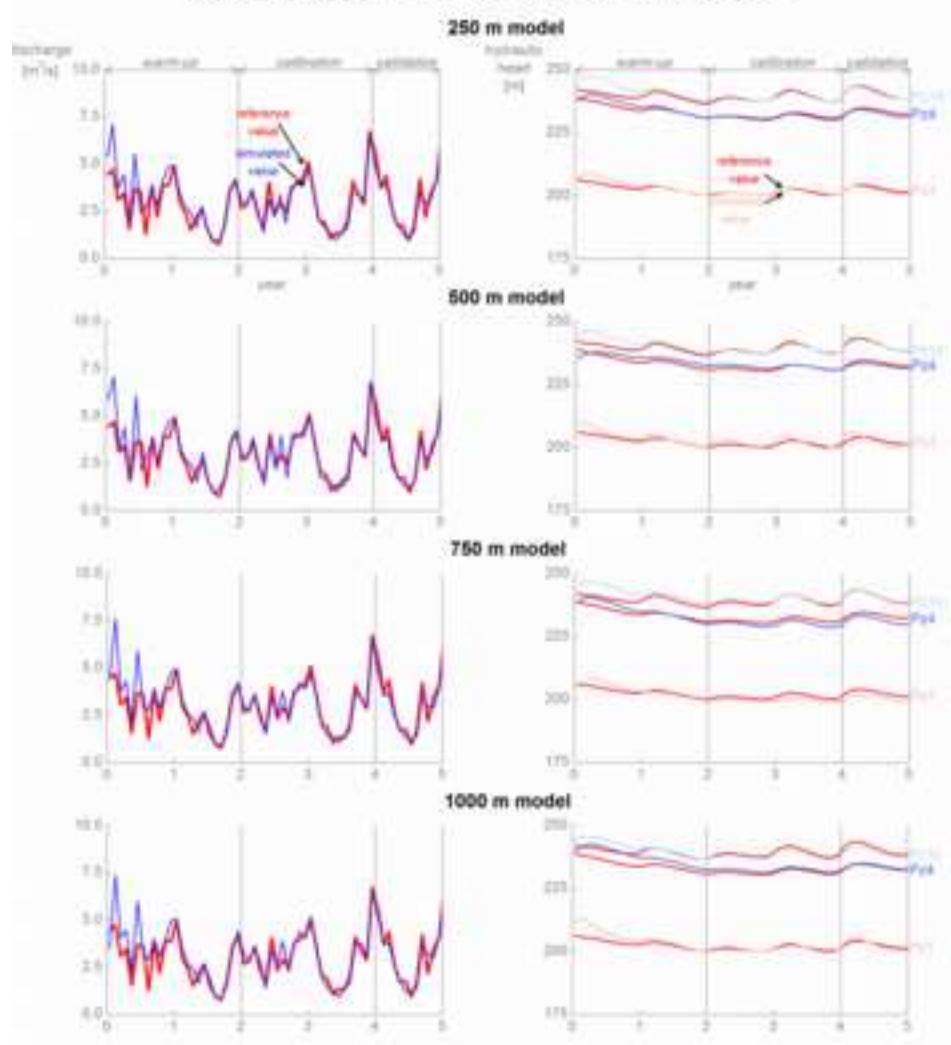
[Click here to download high resolution image](#)



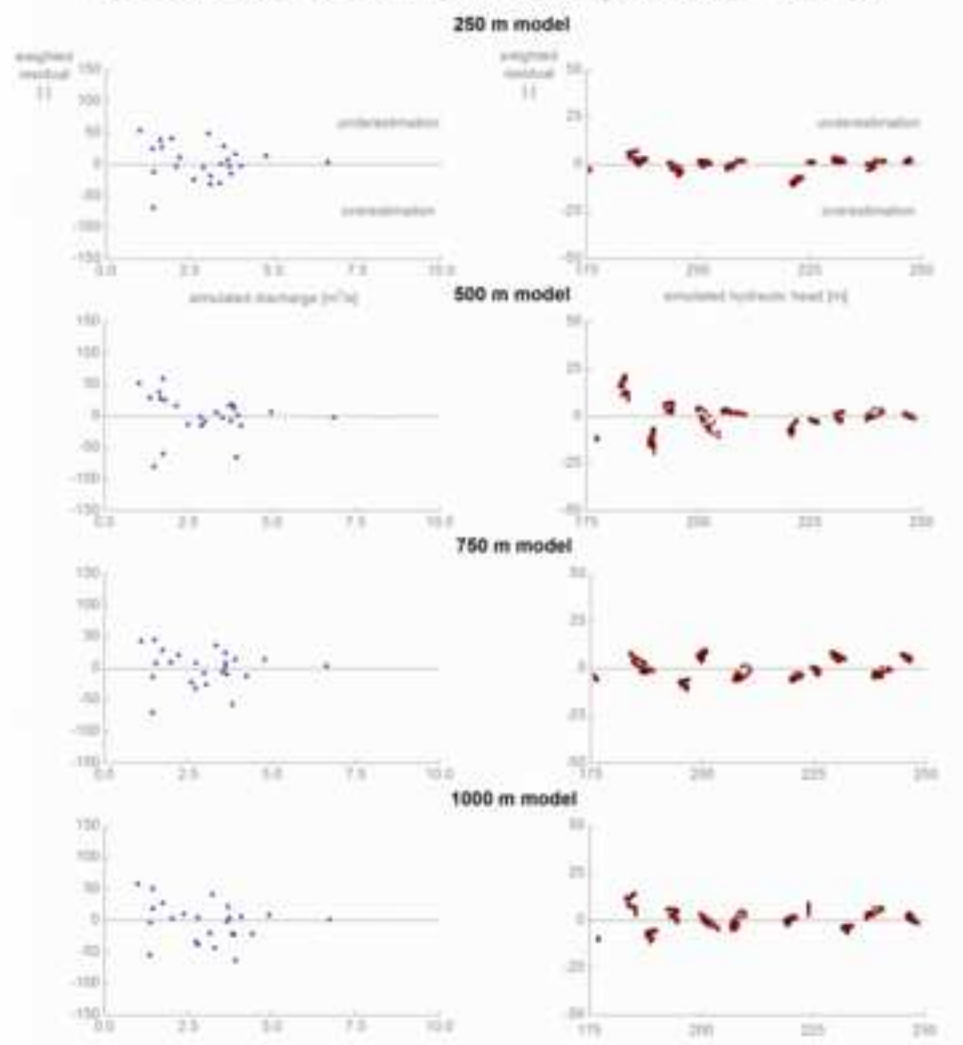
Figure_6

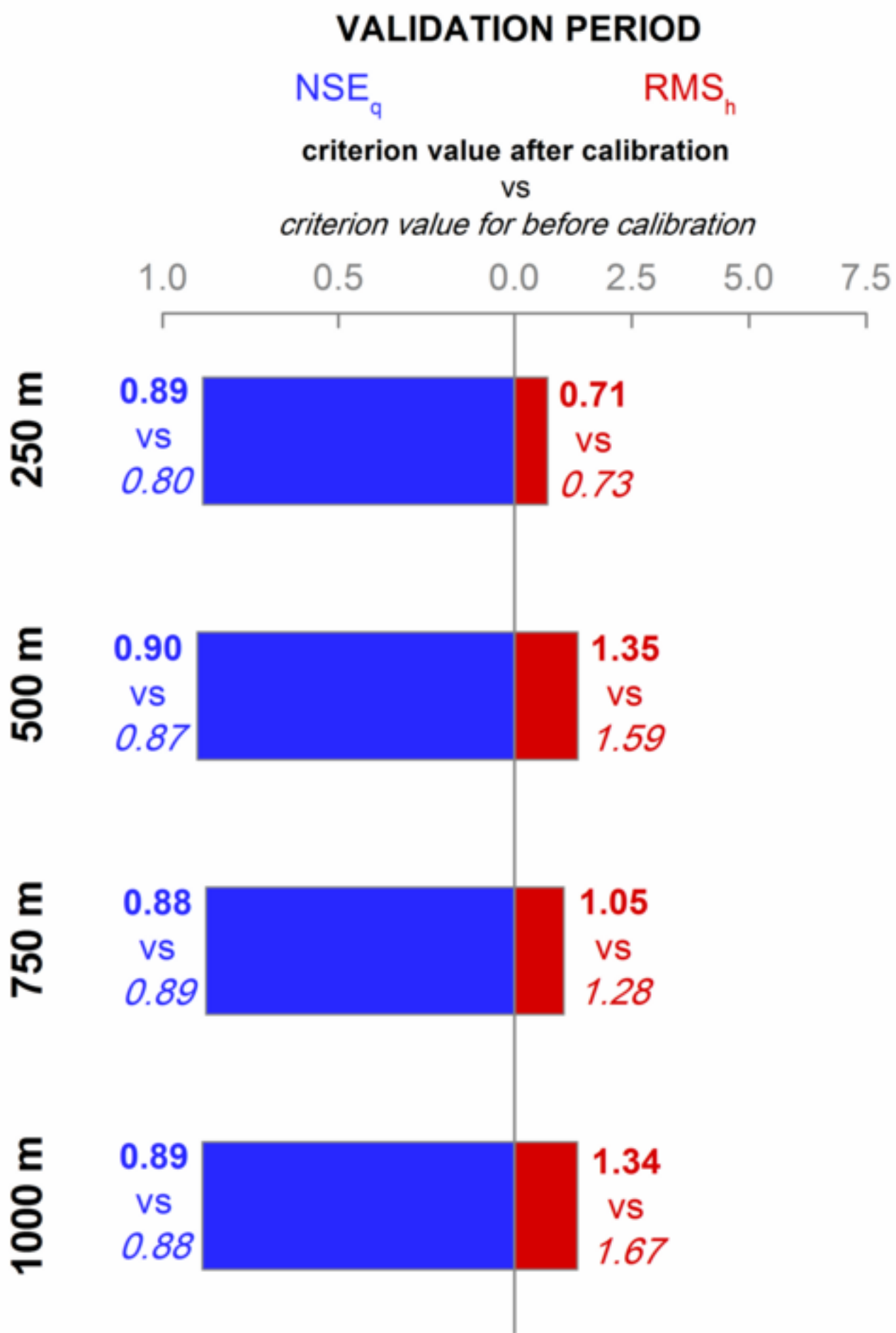
[Click here to download high resolution image](#)

A. reference vs simulated values



B. simulated values vs weighted residuals

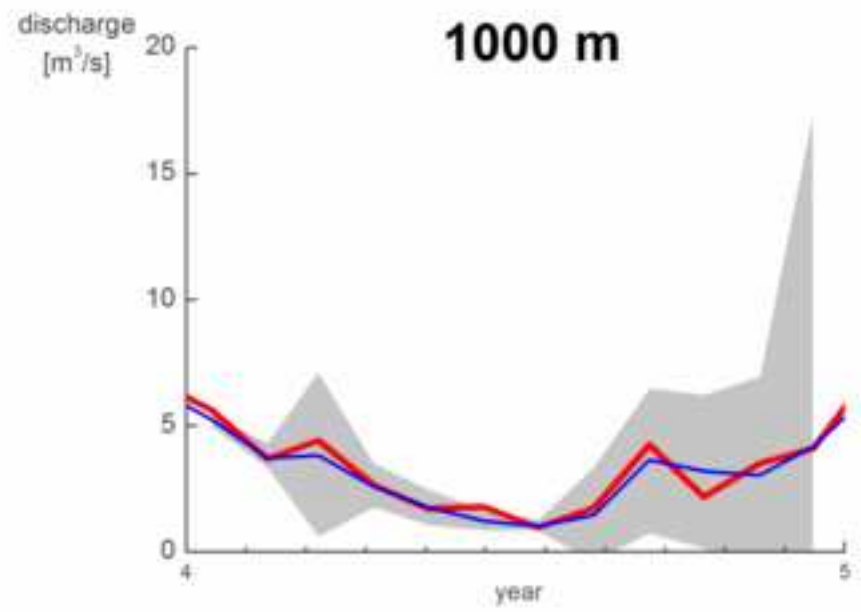
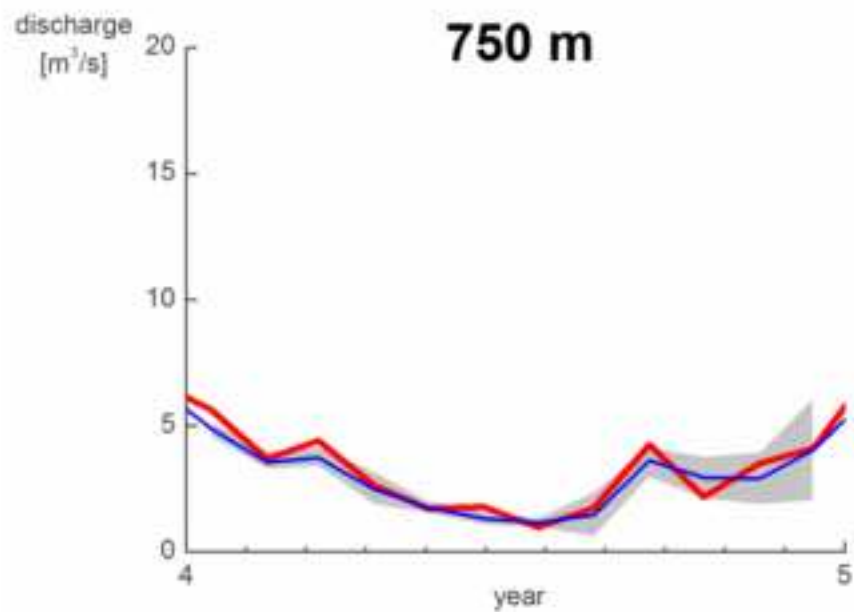
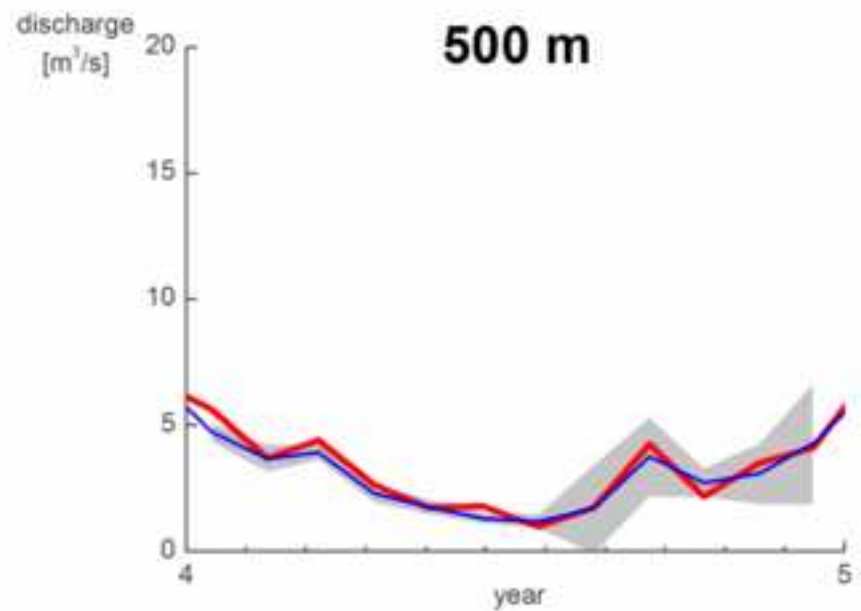
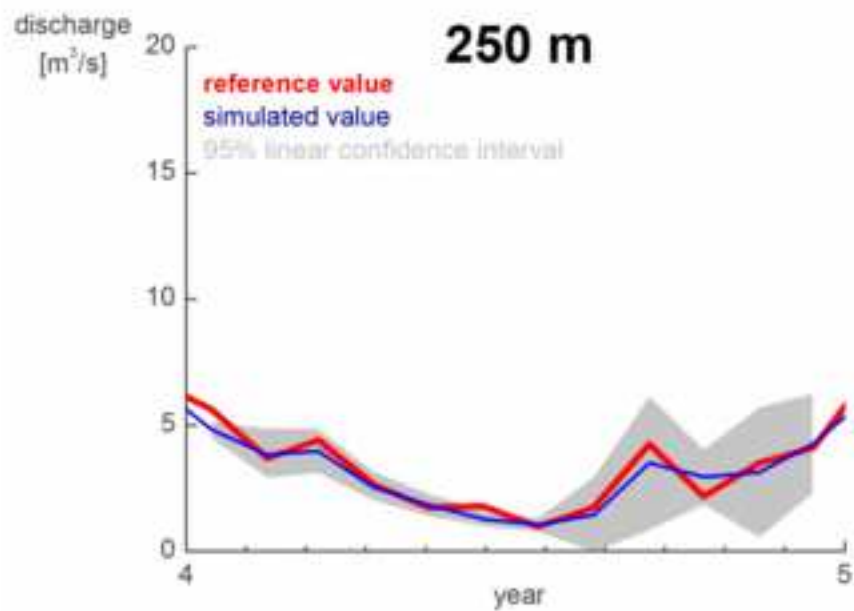




Figure_8

[Click here to download high resolution image](#)

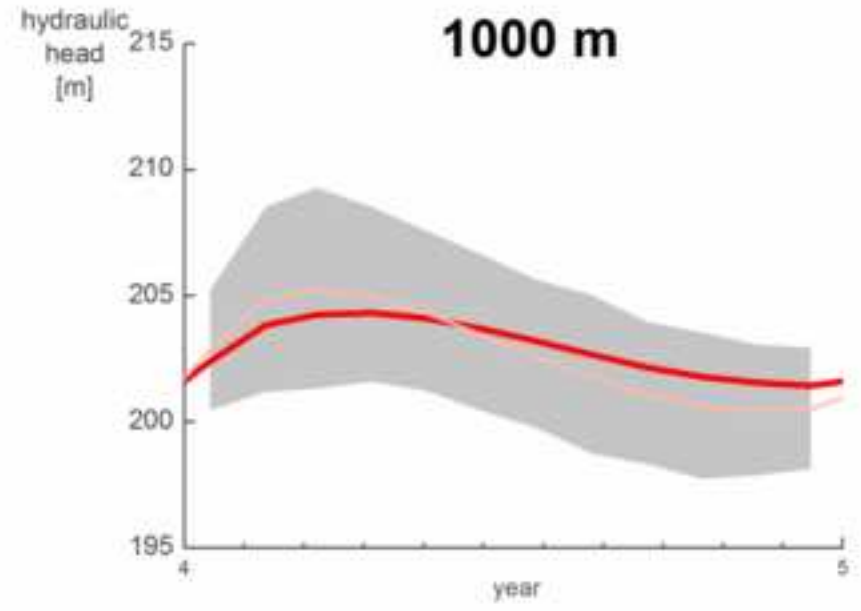
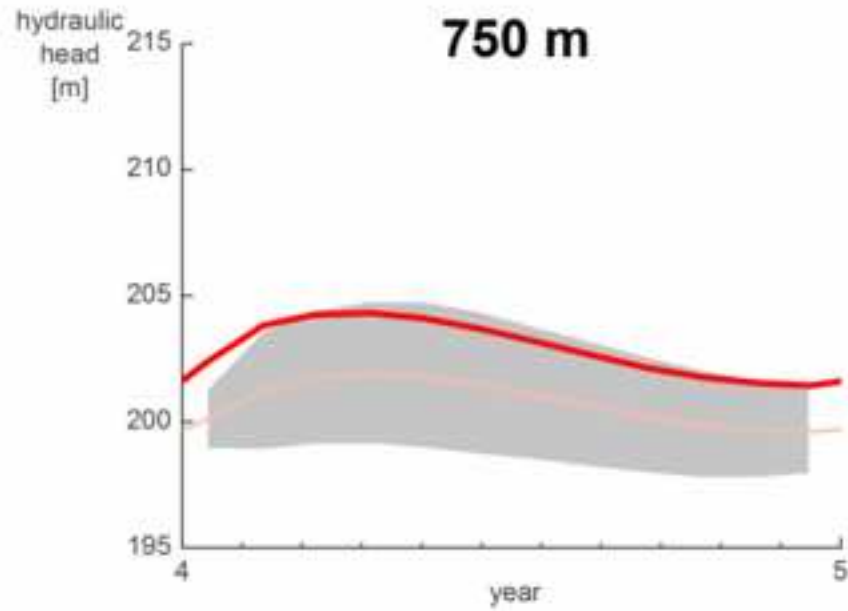
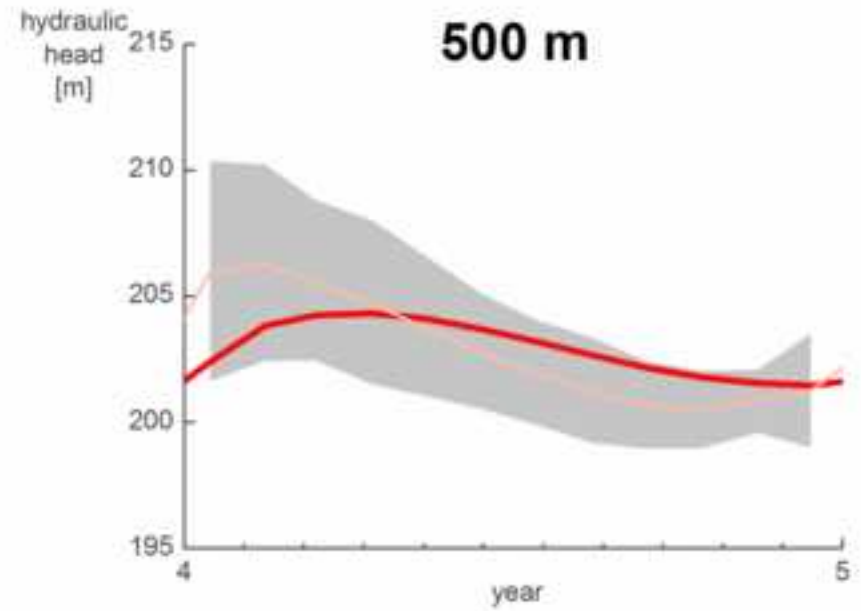
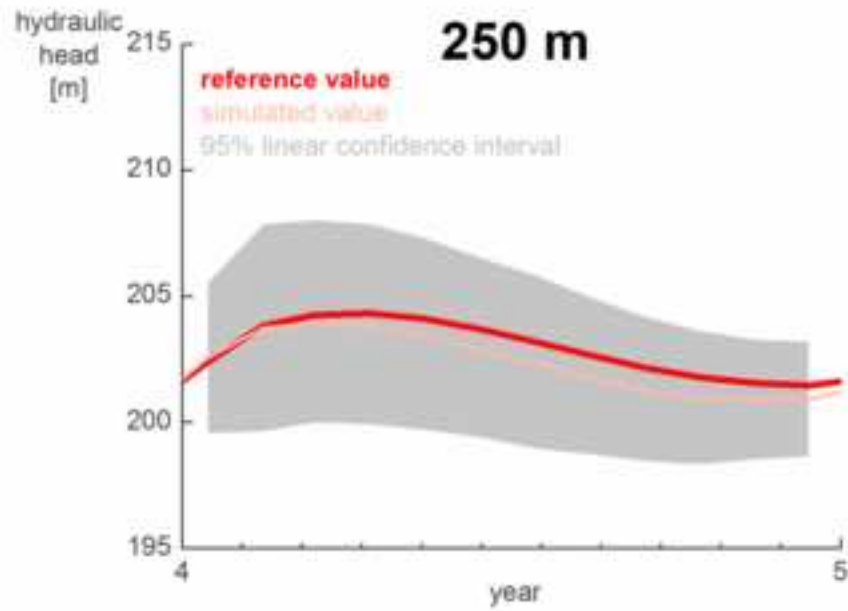
Q_{outlet}



Figure_9

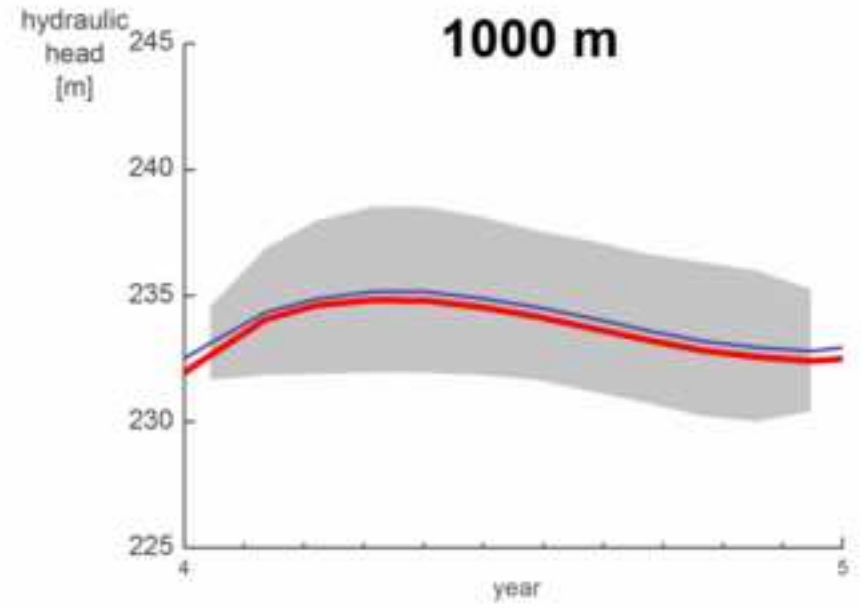
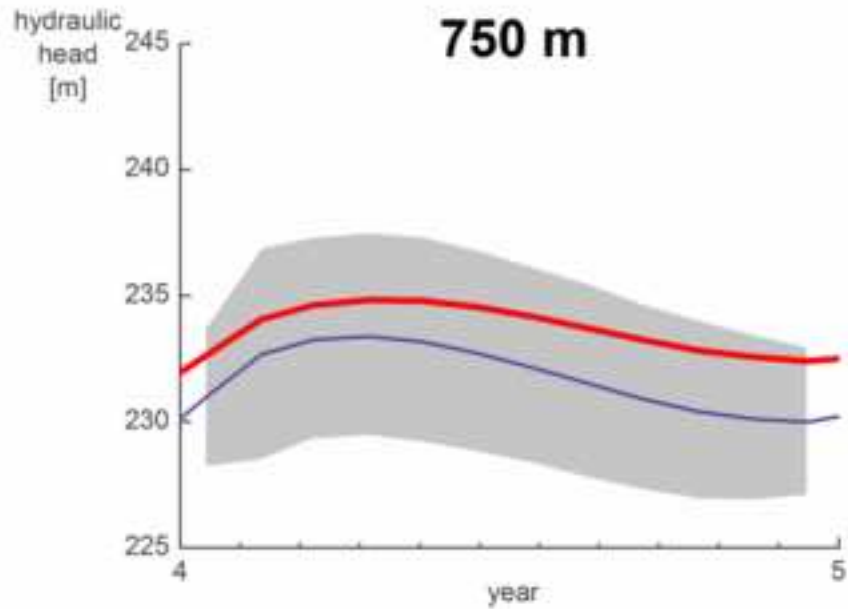
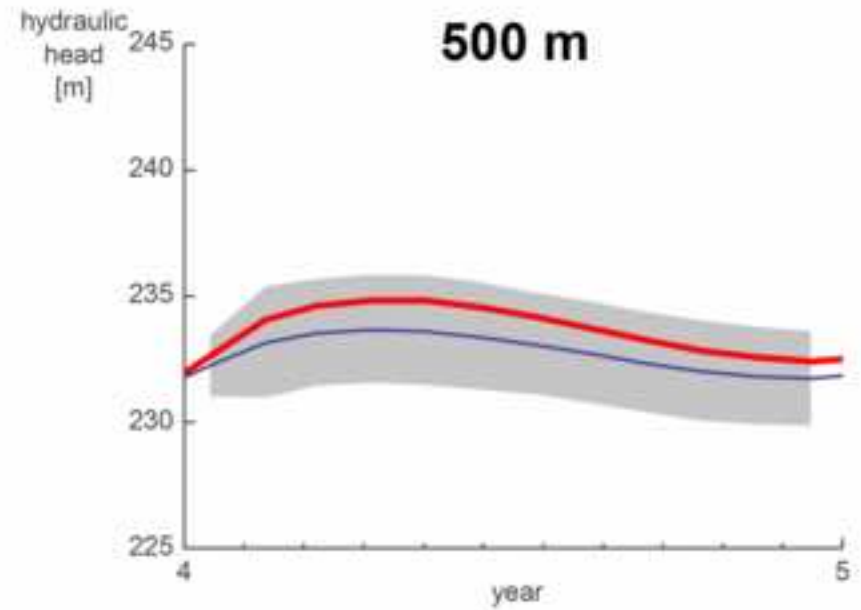
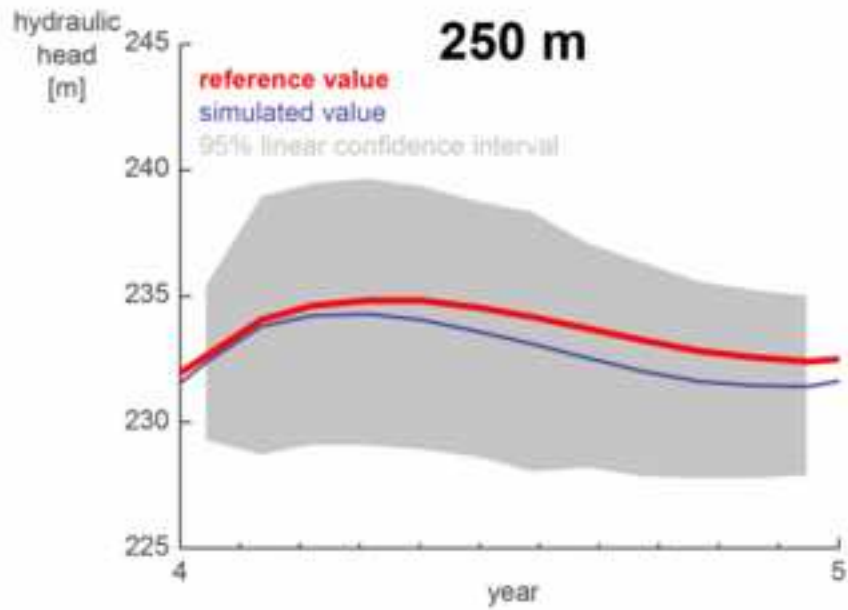
[Click here to download high resolution image](#)

Pz1



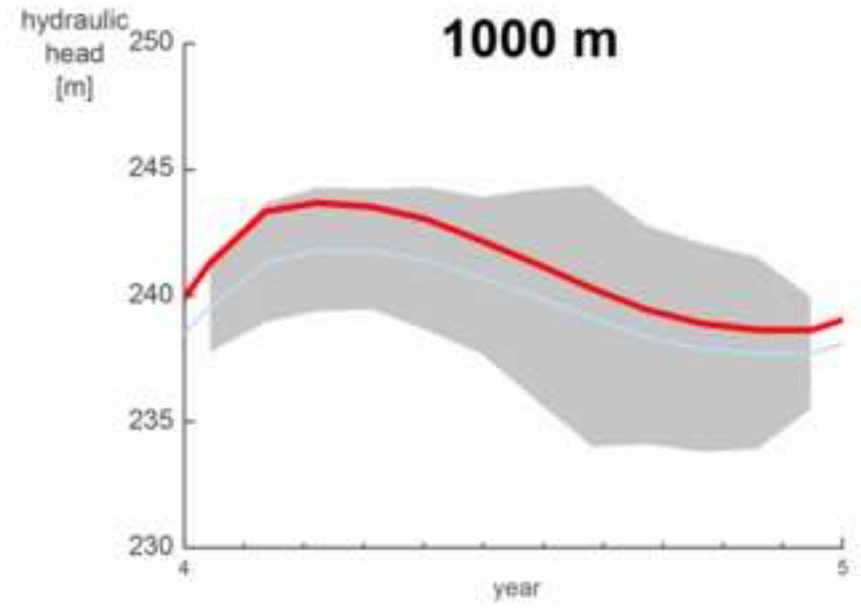
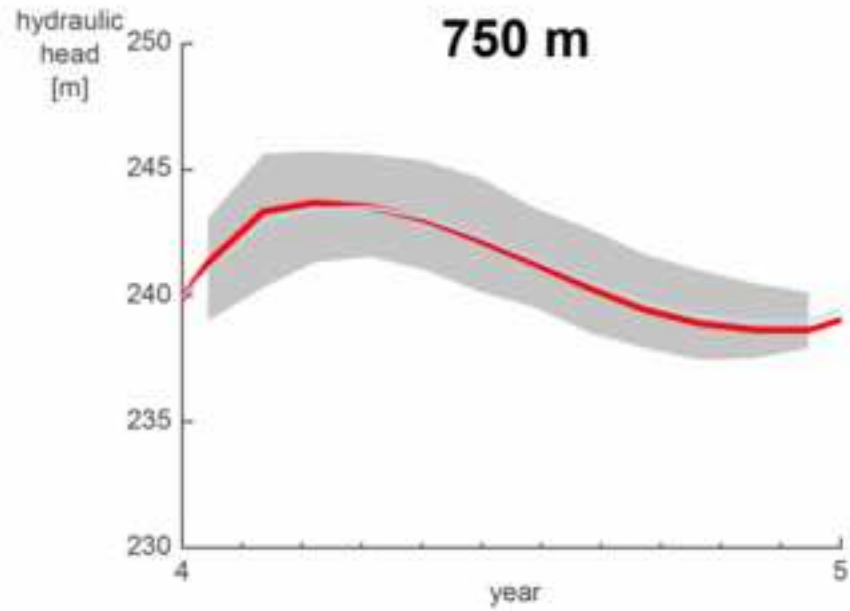
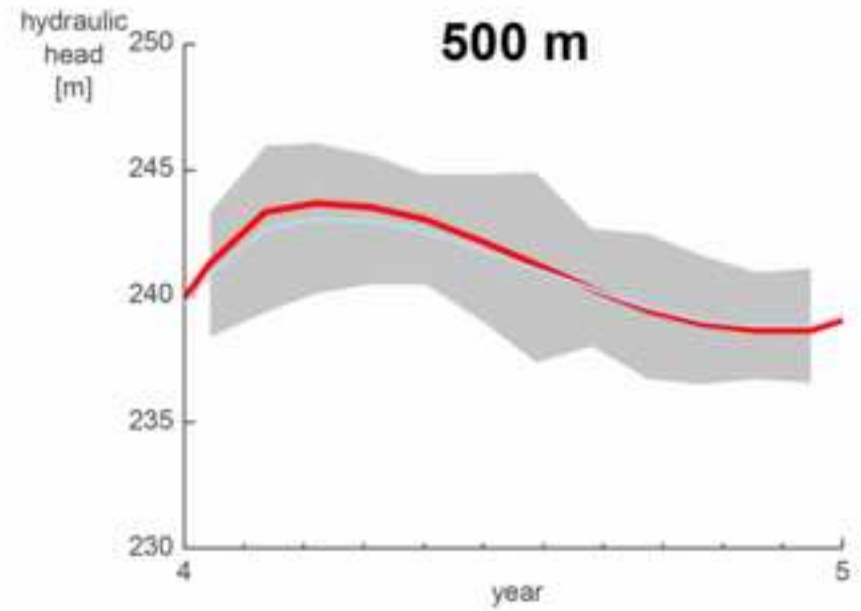
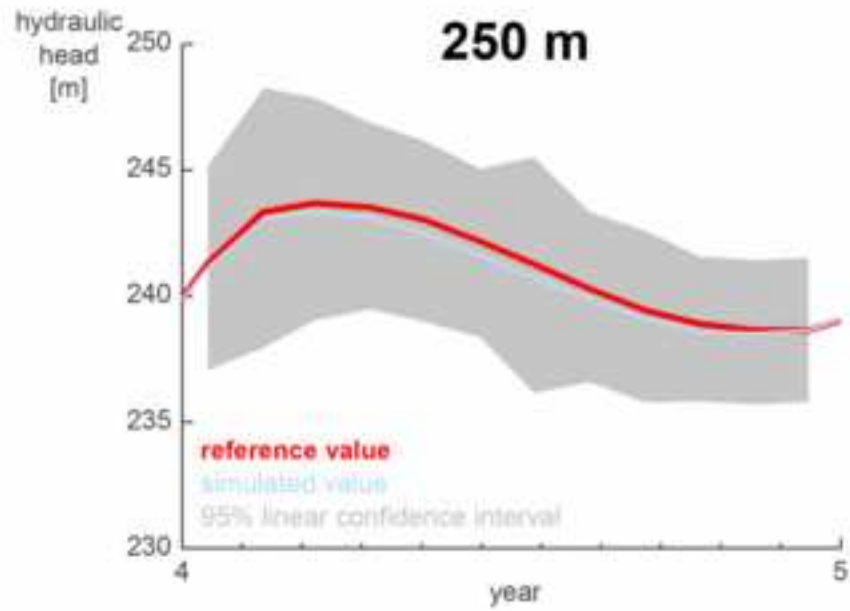
Figure_10
[Click here to download high resolution image](#)

Pz4



Figure_11
[Click here to download high resolution image](#)

Pz10



Table[Click here to download Table: Tables.docx](#)**Table 1**

	Number of nodes	Number of elements	Execution time
250 m	61,884	107,536	6.08 h
500 m	16,200	27,544	0.71 h
750 m	7,245	12,040	0.20 h
1000 m	4,302	7,016	0.13 h

Table 2

	250 m	500 m	750 m	1000 m
[-]	0.77	0.83	0.76	0.72
Gupta's 1st term [-]	2.26	2.33	2.53	2.59
Gupta's 2nd term [-]	1.37	1.47	1.77	1.87
Gupta's 3rd term [-]	0.11	0.03	0.00	0.00
[%]	-14.32	-7.00	-1.32	1.23
[%]	-8.65	-2.01	8.27	11.66
[%]	5.04	11.82	20.27	21.83

Table 3

		250 m		
		min	mean	max
	[%]	-1.03	-0.17	0.50
	[%]	-91.67	-9.40	24.39
	[%]	-89.96	-17.69	3.66
		500 m		
		min	mean	max
	[%]	-1.96	-0.15	1.18
	[%]	-91.67	-4.22	51.35
	[%]	-89.96	-12.62	24.70
		750 m		
		min	mean	max
	[%]	-2.18	-0.36	0.49
	[%]	-91.67	-12.46	19.51
	[%]	-78.26	-19.83	2.76
		1000 m		
		min	mean	Max
	[%]	-2.86	-0.62	0.96
	[%]	-98.04	-21.46	21.62
	[%]	-94.41	-24.97	10.98

Table 4

	250 m	500 m	750 m	1000 m
[-]	0.93	0.93	0.93	0.91
Gupta's 1st term [-]	1.87	2.01	1.87	1.91
Gupta's 2nd term [-]	0.94	1.08	0.94	1.00
Gupta's 3rd term [-]	0.01	0.00	0.00	0.00
[%]	-3.44	-0.28	-1.82	0.26
[%]	-4.81	-2.49	0.28	4.62
[%]	-2.44	1.71	-2.42	-4.25

Table 5

		250 m		
		min	mean	max
	[%]	-0.48	0.06	0.90
	[%]	-95.24	-15.08	49.62
	[%]	-88.46	-7.57	38.43
		500 m		
		min	mean	max
	[%]	-1.09	0.31	1.47
	[%]	-85.71	-7.83	98.20
	[%]	-80.77	-3.78	93.93
		750 m		
		min	mean	max
	[%]	-0.96	-0.04	0.79
	[%]	-36.11	5.14	52.74
	[%]	-28.45	18.92	96.15
		1000 m		
		min	mean	max
	[%]	-1.72	-0.16	0.93
	[%]	-98.04	-21.40	61.26
	[%]	-96.50	-16.89	63.16

Table 6

	reference model	250 m	500 m	750 m	1000 m
K_I [m/s]	5.00×10^{-7}	1.00×10^{-6}	1.66×10^{-6}	1.92×10^{-7}	7.38×10^{-7}
K_{IV} [m/s]	1.00×10^{-4}	9.91×10^{-5}	1.04×10^{-4}	9.75×10^{-5}	8.05×10^{-5}
$\alpha_{VG-IV,V,VI}$ [1/m]	3.65×10^{-2}	4.98×10^{-2}	4.56×10^{-2}	5.04×10^{-2}	3.04×10^{-2}
$\beta_{VG-I,II}$ [-]	1.45	1.47	1.32	1.53	1.39
$\beta_{VG-IV,V,VI}$ [-]	1.83	2.21	2.42	2.18	1.77