# Accepted Manuscript

Calibration and uncertainty analysis of the SWAT model using Genetic Algorithms and Bayesian Model Averaging

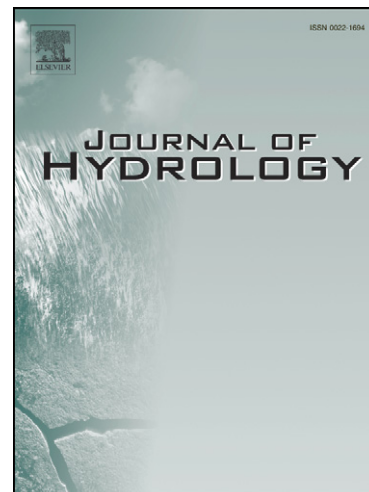Xuesong Zhang, Raghavan Srinivasan, David Bosch

Please cite this article as: Zhang, X., Srinivasan, R., Bosch, D., Calibration and uncertainty analysis of the SWAT model using Genetic Algorithms and Bayesian Model Averaging, *Journal of Hydrology* (2009), doi: 10.1016/j.jhydrol.2009.06.023

1  **Calibration and uncertainty analysis of the SWAT model using**

2  **Genetic Algorithms and Bayesian Model Averaging**

3

4  Xuesong Zhang[1], Raghavan Srinivasan[2], David Bosch[3]

5  [1] Joint Global Change Research Institute, Pacific Northwest National Laboratory, College
6  Park, MD 20740, USA

7  Xuesongzhang2004@gmail.com

8  [2] Spatial Sciences Laboratory, Department of Ecosystem Sciences and Management, Texas
9  A&M University, College Station, TX 77843, USA.

10  [3] Southeast Watershed Research Laboratory, Agricultural Research Service, U.S. Department
11  of Agriculture, Tifton, GA 31793, USA.

12

13  **Abstract** In this paper, the Genetic Algorithms (GA) and Bayesian model averaging (BMA)
14  were used to simultaneously conduct calibration and uncertainty analysis for the Soil and
15  Water Assessment Tool (SWAT). In this combined method, several SWAT models with
16  different structures are first selected; next GA is used to calibrate each model using observed
17  streamflow data; finally, BMA is applied to combine the ensemble predictions and provide
18  uncertainty interval estimation. This method was tested in two contrasting basins, the Little
19  River Experimental Basin in Georgia, USA, and the Yellow River Headwater Basin in China.
20  The results obtained in the two cast studies show that this combined method can provide
21  deterministic predictions better than or comparable to the best calibrated model using GA.
22  66.7% and 90% uncertainty intervals estimated by this method were analyzed. The
23  differences between the percentage of coverage of observations and the corresponding
24  expected coverage percentage are within 10% for both calibration and validation periods in
25  these two test basins. This combined methodology provides a practical and flexible tool to
26  attain reliable deterministic simulation and uncertainty analysis of SWAT.

27  **Key words** optimization; modeling; basin; uncertainty; SWAT

## 1 INTRODUCTION

In recent years, hydrologic models are more and more widely applied by hydrologists and resource managers as a tool to understand and manage ecological and human activities that affect basin systems. Traditionally, the hydrologic models are calibrated to find one optimal hydrologic model with the optimum objective functions (e.g. sum square error). The optimized model is then used to assess water resources practices. The inferences based on a single model implicitly assumes that the probability that the single model generates the data accurately is 1, and neglects the uncertainty inherent in the model selection process (Montgomery and Nyhan, 2008; Raftery and Zheng, 2003). Uncertainty within model output is a major concern, particularly when modeling results are used to set policy. Because of uncertainties associated with input, model structure, parameter, and output, the model predictions are not a certain value, and should be represented with a confidence range (Beven and Binley, 1992, Gupta et al., 1998; Beven, 2000; Beven and Freer, 2001; Beven, 2006; Van Griensven, 2008). Reasonable estimates of prediction uncertainty of hydrologic processes are valuable to water resources and other relevant decision making processes (Liu and Gupta, 2007). Uncertainty estimates are routinely incorporated into Total Maximum Daily Load (TMDL) estimates and are an important part of the TMDL implementation plan (Shirmohammadi et al., 2006). Usually, water management projects are planned and designed using scenarios that fall at the conservative end of the range of plausible outcomes. Over estimation of uncertainty can result in over design of mitigation measures, while under estimation of uncertainty can lead to inadequate preparation for potential situations. In order to successfully apply hydrological models in practical water resources investigations, careful calibration and prediction uncertainty analysis are required (Duan et al., 1992; Beven and Binley, 1992; Vrugt et al., 2003; Yang et al., 2008; Van Griensven et al., 2008).

52    As a physically based hydrologic model that can simulate most of the key hydrologic

53    processes at basin scale, the Soil and Water Assessment Tool (SWAT) (Arnold et al., 1998)

54    has been applied world wide for assessing water resources management (Gassman et al.,

55    2007). In order to efficiently and effectively apply the SWAT model, different calibration and

56    uncertainty analysis methods have been developed and applied to improve the prediction

57    reliability and quantify prediction uncertainty of SWAT simulations (Eckhardt and Arnold,

58    2001; Bekele and Nicklow, 2007; Yang et al., 2007; Harmel and Smith, 2007; Arabi et al.,

59    2007; Kannan et al., 2008). For example, Van Griensven et al. (2003) incorporated the

60    shuffled complex evolution (SCE) algorithm for parameter calibration of SWAT, which was

61    later extended to an uncertainty analysis method known as Sources of Uncertainty Global

62    Assessment using Split SamplES (SUNGLASSES) (Van Griensven et al., 2008). Muleta and

63    Nicklow (2005) combined Genetic Algorithms (GA) and Generalized Likelihood Uncertainty

64    Estimation (GLUE) methods to conduct parameter calibration and uncertainty analysis of

65    SWAT. Yang et al. (2008) compared four uncertainty analysis algorithms, that is GLUE

66    (Beven and Binley, 1992), Sequential Uncertainty Fitting SUFI-2 (Abbaspour et al., 2004),

67    Parameter solutions (ParaSol) (van Griensven and Meixner, 2004), and Markov Chain Monte

68    Carlo (MCMC) based Bayesian analysis techniques for assessing the uncertainty of SWAT

69    predictions. These uncertainty analysis algorithms are differing in philosophy, assumptions,

70    and sampling strategies. Yang et al. (2008) suggested that, if computationally feasible,

71    Bayesian Markov Chain Monte Carlo (MCMC) approaches are most recommendable because

72    of their solid conceptual basis. It is worth noting that the MCMC method requires a large

73    number of SWAT runs. For example, 45,000 runs of SWAT were performed in Yang et al.

74    (2008). Zhang (2008b) test an evolutionary Monte Carlo based MCMC method for SWAT,

75    which took about 200,000 model runs for convergence. Applying the MCMC based methods

76    to assess water resources under future scenarios (e.g. best management practices, and land

77  use/climate change) is very computationally intensive. In the previous uncertainty studies

78  using SWAT, model prediction uncertainty is mainly attributed to parameter values. It is

79  worth noting that the bias and uncertainty result from model structures selection can exert

80  important impact on model prediction (Neuman, 2003; Butts et al., 2004a, 2004b). Butts et al.

81  (2004a) presented an evaluation of model structure on hydrologic modeling uncertainty by

82  selecting different plausible model structures within a general hydrological modeling tool,

83  and emphasize the importance of exploring different model structures as part of the overall

84  modeling approach. The SWAT model provides a hydrologic modeling tool that allows

85  different model structures to be selected for representing different hydrological processes

86  (e.g. potential evapotranspiration, snow routing, and flood routing). The major purpose of this

87  study is to explore ensemble hydrologic simulation and uncertainty analysis using several

88  model structures within the SWAT model framework.

89      Recently, Bayesian Model Averaging (BMA), a method for averaging over different

90  competing models, has been applied to allow incorporating model uncertainty into model

91  prediction. BMA possesses a range of theoretical optimality properties and has shown good

92  performance in reliable prediction and uncertainty analysis in a variety of simulated and real

93  data situations (e.g. weather forecast and hydrologic predictions) (Raftery et al., 2005; Ajami

94  et al., 2006; Duan et al., 2007; Vrugt et al., 2007; Montgomery and Nyhan, 2008). The BMA

95  can be used to examine several competitive models for hydrologic modeling and assessment.

96  In practical applications of SWAT, modelers usually select one or several model structures

97  and choose the best among them. To the best of the authors' knowledge, seldom studies have

98  been conducted to jointly use multiple structures within the SWAT model. In this study, a

99  combined method, which implements the Genetic Algorithms (GA) and BMA, was proposed

100 to conduct calibration and uncertainty analysis of the SWAT model through jointly using

101 multiple model structures. The general procedures for applying GA and BMA to conduct

4

102 ensemble hydrologic predictions applied here are: 1) select the specific model components of

103 SWAT to be examined, here we examined different snow, potential evapotranspiration and

104 flow routing methods; 2) calibrate the parameters for each combination of model components

105 using GA to provide competing models and model results; 3) use BMA to combine the

106 ensemble predictions and provide uncertainty interval estimation. The examination was

107 limited to the snow, potential evapotranspiration and flow routing to present a manageable

108 number of modeling options for illustration purposes. Compared with running thousands of

109 models for assessing management practices or climate / land use change scenarios using

110 MCMC based method, the BMA has the potential to save a large number of runs of SWAT.

111 Two basins were used to test the validity of this framework for providing accurate hydrologic

112 prediction and uncertainty intervals estimation using SWAT. The combination of GA and

113 BMA is expected to provide a practical tool for implementing calibration and uncertainty

114 analysis of computationally intensive hydrologic models.

115 **2. MATERALS AND METHODS**

116 **2.1. Study area description**

117 Two basins, the Little River Experimental Basin (LREB) in the Southeastern USA and

118 Yellow River Headwater Basin (YRHB) in central China were used in this study (Figure 1).

119 The basins were selected to offer a contrast in hydrology for testing purposes. The basic

120 characteristics of the two basins are introduced as follows.

121 The LREB (Figure 1) is the upper 334 km$^2$ of the Little River in Georgia, USA, and is

122 the subject of long-term hydrologic and water quality research by USDA-ARS and

123 cooperators (Sheridan, 1997; Bosch et al., 2007). The LREB is located in the Tifton Upland

124 physiographic region, which is characterized by intensive agriculture in relatively small fields

125 in upland areas and riparian forests along stream channels. The region has low topographic

126 relief and is characterized by broad, flat alluvial floodplains, river terraces, and gently sloping

127 uplands (Sheridan, 1997). Climate in this region is characterized as humid subtropical with an

128 average annual precipitation of about 1167 mm based on data collected by USDA ARS from

129 1971 to 2000. Soils on the basin are predominantly sands and sandy loams with high

130 infiltration rates. Since surface soils are underlain by shallow, relatively impermeable

131 subsurface horizons, deep seepage and recharge to regional ground water systems are

132 impeded (Sheridan, 1997). Land use types include forest (50%), cropland (31%), pasture

133 (10%), water (2%), and urban (7%) (Bosch et al., 2006).

134 The YRHB (Figure 1) is an 114,345 km$^2$ mountainous river basin, which is located in

135 the northeastern part of Tibetan plateau in China. This area is the primary source of water

136 availability for the Yellow River Basin (Liu, 2004). The average elevation is about 4,217 m,

137 and ranges between 2,600 and 6,266 m. The area slopes downward from west to east, ranging

138 from a combined landform of low-mountains and wide valleys with lakes to smooth plateaus.

139 The headwater area has a typical continental alpine cold and dry climate. The annual

140 precipitation amount is around 600 mm and the average annual temperature for the YRHB is

141 near 0°C. In winter the average temperature is below 0°C for most of the weather stations,

142 while in summer the average temperature is above 0°C. This seasonal temperature variation

143 makes snowmelt an important process in this area (Zhang et al., 2008a). This basin is

144 characterized by gently sloping upland, river bed, and swamp and wetland. The major types

145 of soils in this area are clay and loam with relatively low infiltration rates. The major land

146 cover in the study area is grassland, which accounts for approximately 90% of the total area.

147 Other land use/land cover (forest land, rangeland, agriculture land, and bare area) accounts

148 for the remaining 10% of the area.

149 **2.2 SWAT model description**

150 SWAT is a continuous time, physically based hydrological model. SWAT subdivides a

151 basin into sub-basins connected by a stream network, and further delineates Hydrologic

152 Response Units (HRUs) consisting of unique combinations of land cover and soils in each

153 sub-basin. SWAT allows a number of different physical processes to be simulated in a basin.

154 The hydrologic routines within SWAT account for snow fall and melt, vadose zone processes

155 (i.e., infiltration, evaporation, plant uptake, lateral flows, and percolation), and ground water

156 flows. The hydrologic cycle as simulated by SWAT is based on the water balance equation:

157
$$SW_t = SW_0 + \sum_{i=1}^{t}(R_{day} - Q_{surf} - E_a - w_{seep} - Q_{gw}) \tag{1}$$

158 where $SW_t$ is the final soil water content (mm $H_2O$), $SW_0$ is the initial soil water content on

159 day $i$ (mm $H_2O$), $t$ is the time (days), $R_{day}$ is the amount of precipitation on day $i$ (mm $H_2O$),

160 $Q_{surf}$ is the amount of surface runoff on day $i$ (mm $H_2O$), $E_a$ is the amount of

161 evapotranspiration on day $i$ (mm $H_2O$), $w_{seep}$ is the amount of water entering the vadose zone

162 from the soil profile on day $i$ (mm $H_2O$), and $Q_{gw}$ is the amount of return flow on day $i$ (mm

163 $H_2O$). Precipitation in SWAT is divided into rainfall and snowfall. There are three snow

164 routing algorithms available in SWAT, which include the degree day (DD), DD plus

165 elevation band (Fontaine et al., 2002), and the energy balance based SNOW17 models

166 (Zhang et al., 2008a). Surface runoff volume is estimated using a modified version of the Soil

167 Conservation Service (SCS) Curve Number (CN) method (Neitsch et al., 2005a). For

168 evapotranspiration estimation, three options are available in SWAT, that is, Penman-

169 Monteith, Priestley-Taylor, and Hargreaves methods (Neitsch et al., 2005a). A kinematic

170 storage model is used to predict lateral flow, whereas return flow is simulated by creating a

171 shallow aquifer (Arnold et al., 1998). The Variable Storage and Muskingum methods are

172 used for channel flood routing. Outflow from a channel is adjusted for transmission losses,

173 evaporation, diversions, and return flow (Arnold et al., 1998).

174 In the SWAT model, there are numerous parameters to be calibrated to match the

175 simulated and observed flow. Van Liew et al. (2007) tested the suitability of SWAT for the

7

176 Conservation Effects Assessment Project in several USDA Agricultural Research Service

177 basins. In the study conducted by Van Liew et al. (2007), eleven parameters were identified

178 as sensitive for the LREB. These eleven parameters (Table 1) were adjusted by the GA for

179 the LREB in this study. In the YRHB, five parameters (i.e. CN2, ESCO, SURLAG,

180 GW_REVAP, and ALPHA_BF) were adjusted for the calibration according to Zhang et al.

181 (2008a). The general description of the parameters used for the calibration is shown in Table

182 1. The parameters' ranges were limited according to van Griensven et al. (2006) and Neitsch

183 et al. (2005b).

184 **2.3 Genetic Algorithms**

185 Zhang et al. (2009c) compared five global optimization algorithms for parameter

186 calibration of SWAT in four basins, and their results show the advantage of GA over other

187 algorithms for calibrating SWAT. Genetic Algorithms are stochastic search procedures

188 inspired by evolutionary biology of natural selection and genetics (Holland, 1975; Goldberg,

189 1989), such as inheritance, mutation, selection, and crossover. The implementation of GA

190 starts with initializing a population of candidate solutions (called chromosomes) which are

191 randomly sampled from the feasible parameter space. In each generation, the individual

192 chromosomes are selected through a fitness-based process, where the more fit chromosomes

193 in the population are preferred to be selected to reproduce new promising offspring. Next, a

194 new generation population of chromosomes is generated from these selected ones using

195 crossover and mutation operations. The crossover operator chooses "parent" solutions and

196 exchange important building blocks of two parent chromosomes to generate new "offspring"

197 solutions. The "offspring" solutions are then randomly mutated to increase the diversity of

198 new population. Through a steady-state-delete-worst plan (Reca and Martinez, 2006), the

199 fitter chromosomes among the old and new population are input into next generation for

200 evolution. This generational evolution of the parameter solutions is repeated until a maximum

8

201  number of model evaluations are reached. With flexibility and robustness, GAs have been

202  successfully applied to solve complex nonlinear programming problems in many science and

203  engineering branches (Reca and Martinez, 2006). Following Schaffer et al. (1989) and Reca

204  and Martinez (2006), the crossover rate was set to 0.5 and mutation rate was the reciprocal of

205  the parameter dimension. Settings of population size and maximum model runs can

206  substantially affect the performance of GA for calibrating SWAT, a small population size of

207  50 and a maximum number of SWAT runs of 5000 were selected in this study following

208  Zhang et al. (2009c).

209  **2.4 Bayesian Model Averaging**

210  In hydrologic modeling, there are many ensemble based methods that can merge

211  information from multiple sources (e.g. modeling results from different models and observed

212  data from different sources). One simple method is the arithmetic mean method, which

213  simply averages the predictions from several sources equally to obtain the ensemble mean

214  prediction. This method has shown more reliable prediction than single model prediction

215  (Raftery et al. 2005; Hsu et al., 2008). Recently, advanced BMA was proposed to combine

216  multiple weather and hydrologic models results to provide more reliable predictions (e.g.

217  Raftery et al. 2005; Ajami et al., 2006; Duan et al., 2007; Vrugt et al., 2007). BMA is a

218  standard approach to inference in the presence of multiple competing models (Raftery and

219  Zheng, 2003). This approach has been used to infer probabilistic predictions that possess

220  more skill and reliability than the original ensemble members produced by several competing

221  models (Duan et al., 2007). In BMA, the probabilistic distribution of a hydrologic prediction

222  $y$ is the weighted average of the posterior distribution of each model under consideration.

223  Raftery et al. (2005) extended BMA from statistical models to weather forecast models. In

224  the following, the BMA framework developed by Raftery et al. (2005) was introduced. The

225  BMA prediction probability distribution can be represented as

9

226

$$p(y \mid f_1, f_2, \ldots, f_K) = \sum_{k=1}^{K} w_k g(y \mid f_k) \qquad (2)$$

228     where $K$ is the number of competing models and $k$ is the index of each model. $f_k$ denote the

229     bias corrected prediction of a candidate model $M_k$. $w_k$ is $p(f_k \mid D)$, the posterior

230     probability of model prediction $f_k$, also known as the likelihood of model prediction $f_k$

231     being the correct prediction given the observational data, $D$. $w_k$ is nonnegative and with a

232     sum ($\sum_{k=1}^{K} w_k$) of 1. $g(y \mid f_k)$ represents the conditional probability distribution function

233     (PDF) of $y$ conditional on $f_k$. Usually, the conditional distribution $g(y \mid f_k)$ can be

234     represented as a normal distribution, $N(a_k + b_k f_k, \sigma_k^2)$, where $a_k$ and $b_k$ are regression

235     coefficients obtained through least square linear regression. Following Raftery et al. (2005)

236     and Duan et al. (2007), the BMA predictions mean and variance can be calculated as

$$E(y \mid f_1, f_2, \ldots f_K) = \sum_{k=1}^{K} w_k (a_k + b_k f_k) \qquad (3)$$

$$Var(y \mid f_1, f_2, \ldots f_K) = \sum_{k=1}^{K} w_k \left[ (a_k + b_k f_k) - \sum_{i=1}^{K} w_i (a_i + b_i f_i) \right]^2 + \sum_{k=1}^{K} w_k \sigma_k^2 \quad (4)$$

239     where $\sigma_k^2$ is the variance associated with model prediction $f_k$ with respect to calibration data

240     $D$. The BMA prediction mean is the weighted average of individual predictions weighted by

241     the likelihood $p(f_k \mid D)$. It can be viewed as a deterministic prediction and compared with

242     other individual predictions in the ensemble and ensemble mean. The two terms of the right-

243     hand side of equation (4) represent the between-prediction variance and within-prediction

244     variance, respectively. The BMA predicts spread-error correlation, and also accounts for the

245     possibility that ensembles may be underdispersive, which is usually the case in ensemble

246     predictions (Raftery et al., 2005).

247    In order to apply the BMA method, the weights $w_k$ and variance $\sigma_k^2$ need to be

248    estimated. In this study, the maximum likelihood estimation (MLE) method was adopted

249    following Raftery et al. (2005). Let $\theta = \{w_1, w_2, \ldots, w_K, \sigma_1^2, \sigma_2^2, \ldots, \sigma_K^2\}$. The log form of the

250    likelihood needs to be maximized is

251
$$\ell(\theta) = \log\left[\sum_{k=1}^{K} w_k g(y \mid f_k)\right] \tag{5}$$

252    It is difficult to analytically maximize this log likelihood. In this study, the Expectation and

253    Maximization (EM) was used to find the maximum likelihood estimator. EM algorithm is

254    iterative. It starts with a initial guess of $\theta^0$. Then the EM algorithm alternates between the

255    Expectation step and Maximization step to update the estimation of $\theta^{Iter}$, where $Iter$ is the

256    iteration number. Finally, the Expectation step and Maximization step converge and are

257    stopped when there is no significant change, measured by a small tolerance value, between

258    two consecutive iterative log likelihood estimations. Following Raftery et al. (2005) and

259    Duan et al. (2007), the procedures of applying EM algorithm to estimate $w_k$ and $\sigma_k^2$ are

260    briefly described in Appendix A. The probabilistic predictions of the variable of interest can

261    be derived based each individual deterministic prediction and its weight and variance. The

262    procedures used in this study to generate probabilistic predictions at each time step $t$ are

263    briefly described as follows (Gelman et al., 2003): i) select an individual competing model

264    ($M_k$) with the probability proportional to its weight; ii) draw a replication $y^{rep}$ from

265    $g(y_t \mid f_{k,t})$; iii) repeat steps i and ii to obtain 1000 values that represent the distribution of

266    $y_t$, with which the uncertainty intervals can be derived. For example, the 90% interval is

267    within the range of the 5% and 95% quartiles. Similarly, other uncertainty intervals with

268    different expected coverage percentage can be derived straightforward.

269    **2.5 Generating competing hydrologic predictions of SWAT**

270    Hydrologic environments are open and complex, rendering them prone to multiple

271    interpretations and mathematical descriptions (Neuman, 2003). In practical application of

272    hydrologic model, modelers typically select a single model among the several choices that is

273    assumed to best represent the hydrologic system. The major advantage of BMA is to jointly

274    use several model structures identified as plausible by the modelers. For the selection of

275    candidate models for BMA, it is suggested to use previous research and theory to specify the

276    set of model structures that are plausible and supported by data (Gelman and Rubin, 1995;

277    Raftery et al., 2005; Duan et al., 2007; Vrugt et al., 2007; Montgomery and Nyhan, 2008). In

278    this study, we followed the methodology used in previous literature on model structures

279    selection. In the selection of model structure, we used the information provided in previous

280    literature on SWAT (Neitsch et al., 2005a, 2005b) and the actual watershed characteristics.

281    The purpose of this paper is to illustrate the application of GA and BMA for combining

282    several plausible model structures within SWAT framework. It is out of the scope of this

283    study to explore all possible model structures.

284    The SWAT model has several options for setting its model structures. Different

285    evapotranspiration, snow accumulation and melt, and flow routing algorithms are available

286    within the SWAT model system. In the LREB, as snowfall and melt is not an important

287    process, we set up SWAT model structures by selecting different evapotranspiration and flow

288    routing algorithms. For the potential evapotranspiration, Penman-Monteith "PM", Priestley-

289    Taylor "PT", and Hargreaves "HG" were selected. For flow routing, Variable Storage "VS"

290    and Muskingum "MK" were selected. Thus, SWAT_PM_VS denotes SWAT with Penman-

291    Monteith potential evapotranspiration estimation and variable storage flow routing. A total of

292    six   model   structures   were   defined,   that   is,   SWAT_PM_VS,   SWAT_PM_MK,

293    SWAT_PT_VS, SWAT_PT_MK, SWAT_HG_VS, and SWAT_HG_MK. The evaluation

294    time scale selected for the LREB was day. In the YRHB, we only choose three models with

12

295    different snowfall and melt algorithms, because snow processes are significant (Zhang et al.,

296    2008a) and the evaluation time scale was month. Previous studies (e.g. Fontaine et al., 2002;

297    Zhang et al., 2008a) have shown that the SWAT model simulation is sensitive to snow

298    routing methods in mountainous basin. The snow routing methods used in this study include

299    the degree day "DD", DD plus elevation band "ELEV", and the energy based SNOW17

300    methods. The SWAT models with different snow modules are represented as SWAT-DD,

301    SWAT-ELEV, and SWAT-SNOW17, respectively. The GA was applied to optimize the

302    SWAT models with different structures in the LREB and YRHB. In the LREB, daily

303    streamflow from 1999 to 2000 was used to calibrate model and daily streamflow from 2001

304    to 2002 was used to validate the model.  Watershed weighted annual precipitation for this

305    period for LREB varied from a high of 1085 mm observed in 2000 to a low of 884 mm

306    observed in 1999.  Precipitation and flow for both the calibration and validation periods were

307    slightly below long term means. For the YRHB, monthly flow from 1976 to 1985 was used to

308    calibrate model and monthly flow from 1986 to 1990 was used to validate the model.

309    Precipitation for this period varied from 653 mm to 482 mm. Precipitation and flow of the

310    selected periods in the YRHB are very close to long term average conditions. The calibrated

311    models with smallest sum square error in the LREB and YRHB serve as competing models

312    for the BMA, and the BMA mean and prediction uncertainty interval are calculated.

313    **2.5 Statistical criteria for evaluating the performance of hydrologic prediction**

314    Different statistical criteria were used to evaluate the individual SWAT model

315    predictions, ensemble mean, BMA mean, and the uncertainty intervals obtained by the BMA.

316    Following Santhi et al. (2001) and Moriasi et al. (2007), the evaluation coefficient for

317    deterministic predictions include Percent Bias ($PBIAS$), Coefficient of Determination ($R^2$),

318    and Nash-Sutcliffe Efficiency ($NSE$). $PBIAS$ is calculated as

319
$$PBIAS = \left(\sum_{t=1}^{T}(f_t - y_t)\right) \bigg/ \sum_{t=1}^{T} y_t \times 100 \qquad (6)$$

320 where $f_t$ is the model simulated value at time unit $t$, $y_t$ is the observed data value at time

321 unit $t$, and $t = 1,2,\ldots,T$. *PBIAS* measures the average tendency of the simulated data to be

322 larger or smaller than their observed counterparts (Gupta et al., 1999). *PBIAS* values with

323 small magnitude are preferred. Positive values indicate model overestimation bias, and

324 negative values indicate underestimation model bias (Gupta et al., 1999).

325 The formula for calculating coefficient $R^2$ is

326
$$R^2 = \left\{ \sum_{t=1}^{T}(y_t - \bar{y})(f_t - \bar{f}) \bigg/ \left[\sum_{t=1}^{T}(y_t - \bar{y})^2\right]^{0.5}\left[\sum_{t=1}^{T}(f_t - \bar{f})^2\right]^{0.5} \right\}^2 \qquad (7)$$

327 where $\bar{y}$ is the mean of observed data value for the entire time period of the evaluation, $\bar{f}$ is

328 the mean of simulated data value for the entire time period of the evaluation. The other

329 symbols have the same meaning defined above. $R^2$ is equal to the square of the Pearson's

330 product-moment correlation coefficient (Legates and McCabe, 1999). It represents the

331 proportion of the total variance in the observed data that can be explained by the model. $R^2$

332 ranges between 0.0 and 1.0. Higher values mean better performance.

333 *NSE* is calculated as

334
$$NSE = 1.0 - \sum_{t=1}^{T}(y_t - f_t)^2 \bigg/ \sum_{t=1}^{T}(y_t - \bar{y})^2 \qquad (8)$$

335 *NSE* indicates how well the plot of the observed value versus the simulated value fits the 1:1

336 line, and ranges from $-\infty$ to 1 (Nash and Sutcliffe, 1970). The larger the *NSE* values, the

337 better model performance.

338 In hydrologic modeling, different types of uncertainty limits can be recognized (e.g.

339 Beven, 2006; Liu and Gupta, 2007). In this study, we are concerned with the modeling

340  uncertainty and predictive uncertainty (Liu and Gupta, 2007). The modeling uncertainty

341  limits, obtained through calibrating hydrologic models to match observed streamflow data,

342  were expected to include a specified proportion of the calibration data set. The predictive

343  uncertainty limits, obtained through applying the calibrated models to another independent

344  data set, were expected to contain a specified proportion of future observations. In this study,

345  the percentage of coverage (POC) of observations in the uncertainty interval was used to

346  evaluate the uncertainty intervals obtained by the BMA scheme. The smaller difference

347  between POC and the expected coverage percentage of an uncertainty interval indicate better

348  performance of the estimated uncertainty interval. For a 90% uncertainty interval, which is

349  expected to include 90% of the observed data, the POC value closer to 90% indicate the

350  better performance of the uncertainty interval estimation.

351  **3. RESULTS AND DISCUSSION**

352  **3.21 Calibration and uncertainty analysis results in the LREB**

353       The evaluation coefficients of the simulated daily streamflow by different prediction

354  techniques in the LREB are listed in Table 2. The two sample Kolmogorov–Smirnov test (K–

355  S test) (Chakravarti et al., 1967) reveals that the difference between the simulated results by

356  models with default input and those calibrated by GA is significant at a significant level of

357  0.05. This indicates that model calibration can substantially improve model simulation. The

358  calibrated parameters for the six models in the LREB are shown in Table 3, which clearly

359  show that different model structure prefer different parameter values. For example, the

360  calibrated values of CN range between -17% and 20%. For illustration purpose, the simulated

361  daily streamflow by the different methods in March, 1999 and in March, 2001 are shown in

362  Figures 2 and 3 for calibration and validation periods, respectively. The ensemble mean and

363  BMA mean predictions were also plotted for comparison purpose. From Figures 2 and 3,

364  there is obvious difference between the hydrographs simulated by different models,

15

365    especially in the validation period. At a significant level of 0.05, the K-S test results show

366    that there is significant difference between different model simulation results. The evaluation

367    coefficients in Table 2 confirmed the difference between different models. For example, in

368    calibration period, SWAT-HG-VS obtained *PBIAS* of -0.72%, while the *PBIAS* value of

369    SWAT-PM-VS was 24.9%. The performance of calibrated models in validation period is

370    different from that in calibration period. For example, the *PBIAS* values of SWAT-PT-VS

371    increased from 22.94% in calibration period to 46.82% in validation period. Analysis of other

372    evaluation coefficients also shows the difference between model performance in calibration

373    and validation period (Table 2). The difference between model performance in calibration

374    and validation periods is because the hydrologic conditions in validation period may change

375    and do not look exactly like the hydrologic conditions during the calibration period (e.g.

376    Beven, 2006; Liu and Gupta, 2007; Zhang et al., 2009a). The different properties exhibited

377    by various models were combined by the arithmetic ensemble mean and Bayesian model

378    averaging methods. The comparison of the evaluation coefficients of each single model and

379    those of the ensemble based methods indicate the obvious superiority of applying ensemble

380    based methods. Compared with single models predictions, the simple arithmetic ensemble

381    mean obtained better results in terms of $R^2$, and N$SE$ during both calibration and validation

382    period. The BMA outperformed all the other seven methods in terms of all four evaluation

383    coefficients in both calibration and validation periods. The above analysis clearly illustrates

384    the advantage of using ensemble based methods to obtain reliable deterministic streamflow

385    simulation, especially the Bayesian Model Averaging.

386    The 66.7% and 90% uncertainty intervals estimated by the BMA are shown in Figures 4

387    and 5 for calibration and validation periods, respectively. The estimated 66.7% and 90%

388    uncertainty intervals cover about 76.04% and 91.14% of the observed data, respectively, in

389    calibration period, and about 74.41% and 96.53% of the observed data, respectively, in

390  validation period. The absolute differences between the POCs values computed from the

391  uncertainty intervals estimated by the BMA and expected coverage percentages are within

392  10%. In general, the POC values estimated by BMA are matching well with the expected

393  coverage percentage.

394  **3.2 Calibration and uncertainty analysis results in the YRHB**

395  The evaluation coefficients of the simulated monthly streamflow by different prediction

396  techniques in the YRHB are listed in Table 4 for different prediction techniques. The K-S test

397  results indicate that the difference between the simulated results by models with default input

398  and those calibrated by GA is significant at a significant level of 0.05, which emphasize the

399  importance of parameter calibration. The calibrated parameters (Table 5) for the three models

400  also exhibit very different values in the YRHB. Using different snow routing methods can

401  lead to variation of calibrated CN values from 2% to 14%. For illustration purpose, the

402  simulated monthly streamflow by the different methods in 1976 and in 1986 are shown in

403  Figures 6 and 7 for calibration and validation periods, respectively. Similar to the case in the

404  LREB, the hydrographs simulated by the three models with different snow routing algorithms

405  have pronounced differences. The SWAT-DD model consistently underestimates the

406  streamflow, with *PBIAS* values of -17.71% and -17.98% for calibration and validation

407  periods, respectively. The SWAT-SNOW17 model obtained positive *PBIAS* values less than

408  10% for both calibration and validation periods. The arithmetic ensemble mean and BMA

409  mean predictions consistently obtained better performance in terms of $R^2$, and *NSE* than

410  single model based predictions. In terms of *PBIAS*, BMA mean outperformed all the other

411  methods in calibration period, while it performed less than SWAT-ELEV in validation

412  period. But BMA mean still predicted small *PBIAS* value (less than 5%) in the validation

413  period. In the YRHB test case, BMA provided better deterministic prediction than the best

414  ensemble number in calibration period and similar results in validation period.

17

415    The uncertainty intervals with different expected coverage percentages estimated by

416    BMA are shown in Figures 8 and 9 for the calibration and validation periods, respectively.

417    The differences between the estimated POC values by BMA and the corresponding expected

418    coverage percentages are within 6% for both calibration and validation periods. The

419    estimated 66.7% and 90% uncertainty intervals cover about 64.2% and 87.5% of the observed

420    data, respectively, in calibration period, and about 68.67% and 91.67% of the observed data,

421    respectively, in validation period. This good match indicates the validity of using only three

422    ensemble members to estimate the uncertainty of hydrologic predictions.

423    **3.4 Discussion**

424    The test results in the two contrasting basins indicate the combination of GA and BMA

425    holds promise to be an efficient and effective technique to calibrate SWAT model and

426    provide reasonable estimation of prediction uncertainty. The numbers of model runs of

427    SWAT are 30000 and 15000 in LREB and YRHB, respectively. These numbers of model

428    runs reported in this study is much less than those reported in previous studies. For example,

429    two previous studies that applied MCMC for SWAT reported 45000 (Yang et al., 2008) and

430    200000 (Zhang, 2008b) model runs. In addition, in contrast to MCMC methods which usually

431    require thousands of SWAT runs, one only needs to run several competing SWAT models

432    with different model structures to assess water resources effect of different management and

433    global change scenarios. For the computationally intensive SWAT model, the method used in

434    this study has the potential to save enormous computational resources and time. It is still

435    important to note that the time consumed by calibrating one model structure is intensive. We

436    calibrated the candidates SWAT models on a computer with Pentium IV 3 GHZ and 1GB

437    RAM. In the LREB, calibration of each of six model structures took about 3 days. A total of

438    18 days were spent on model calibration for the six model structures in the LREB. In the

439    YRHB, calibration time consumed by SWAT-DD, SWAT-ELEV, and SWAT-SNOW17 was

440　3 days, 5 days, and 25 days, respectively. Given the enormous time consumed by

441　constructing candidate model structures for BMA, using as small number of candidates as

442　possible is very important. We tested the effect of reducing number of models on BMA

443　prediction. In the LREB, we eliminated the candidate model with less *NSE* in sequence until

444　there were only two models remaining. The calculated *PBIAS*, $R^2$, *NSE*, and *POC* values for

445　each combination of model structures are listed in Table 5. The difference between these

446　evaluation coefficients is very small. For example the *NSE* values range between 0.8 and 0.81

447　in calibration period and between 0.84 and 0.86 in validation period. The difference between

448　*POC* values are less than 5% for both 66.7% and 90% intervals. It is worth noting that the

449　*PBIAS* value reached 10% in validation period when using 2 candidate models. This

450　compares to the *PBIAS* values less than 5% for the other combinations of candidate models.

451　Further test in the YRHB show that the evaluation coefficients obtained with two candidate

452　models (SWAT-ELEV and SWAT-SNOW17) are also very close to those calculated using all

453　three models. Overall, reducing number candidate models does not have substantial effect on

454　the performance of BMA in the two case studies. This result is similar to that in Raftery et al.

455　(2005). Considering the relatively large *PBIAS* value obtained by using two candidate models

456　in LREB, it is suggested that three or more model structures are needed for BMA. As

457　hydrologic conditions are varying from site to site, much care should be taken when transfer

458　the results to other basins.

459　　There are several limitations of the method used in this study. It is also worth noting

460　that the BMA mean prediction can not always outperform the other models predictions for all

461　the evaluation coefficients. For example, in the YRHB, the BMA mean predicted larger

462　PBIAS than SWAT-ELEV and performed almost the same as the simple arithmetic ensemble

463　mean in validation period. The K–S test results show that the BMA mean prediction is

464　significantly different from all ensemble members in LREB at a significance level of 0.05.

19

465    While in YRHB, the BMA mean is significantly different from all ensemble members at

466    significance level of 0.2. As significance level of 0.05 is commonly used in hydrologic

467    modeling, the results indicate that the relatively complex BMA analysis did not necessarily

468    show significant improvement. The discrepancy between POC values obtained by the BMA

469    and the expected coverage percentage, which reached about 10% and 6% respectively in the

470    LREB and YRHB, respectively, also shows the BMA methods can be further improved.

471    These inadequacies of the BMA method may be caused by several reasons: i) the uncertainty

472    associated with the input data was not explicitly accounted for. For example, the precipitation

473    uncertainty may have important effect on uncertainty estimation (Kavetski et al., 2006);  ii)

474    the residuals between simulated and observed streamflow data are assumed to independent,

475    which may not be true in real world problems (Kuczera and Parent, 1998; Yang et al., 2007);

476    iii) the prior knowledge of different uncertainty sources, which may affect the uncertainty

477    estimation (Zhang et al., 2009a), was not explicitly considered in the BMA scheme. In the

478    future, incorporating more sources of uncertainty into account (Kuzera et al., 2006) may

479    improve the performance of this method. Methods on incorporating input data uncertainty,

480    obtaining prior knowledge of model, and considering correlation between residuals deserve

481    further research for improving the reliability of SWAT predictions. Another limitation of this

482    method is that the application of GA for parameter estimation took very long time. The

483    expensive computational cost is limiting the use of this method. In the future, incorporating

484    surrogate model (e.g. Zhang et al., 2009b) and parallel computing techniques (e.g. Vrugt et

485    al., 2006) into the model calibration process is a promising research topic.

486    For water resources investigations essential for relevant decision making processes, the

487    predictive uncertainty estimation associated with hydrologic simulation is valuable.

488    Predictive uncertainty limits are dependent on and different from modeling uncertainty. This

489    is because when the calibrated hydrological models are applied to another set of data

490   independent of the calibration data, the hydrologic conditions may change and therefore

491   impact the predictive interval estimation (Beven, 2006; Liu et al., 2008; Zhang et al,. 2009a).

492   The results obtained in the two test basins show that the percentage of coverage values of

493   modeling and predictive uncertainty intervals can be different from each other. In the YRHB,

494   the predictive uncertainty interval included more observed data than the modeling uncertainty

495   intervals. For example, POC value of the 90% interval is 4% less in calibration period than

496   that in validation period. In the LREB, the modeling uncertainty intervals in calibration

497   period included more observed data for 66.7% interval than the corresponding predictive

498   uncertainty intervals in validation period, while the 90% modeling uncertainty interval

499   included about 5% less observed data than the 90% predictive uncertainty interval. Because

500   of the future uncertainties due to natural and anthropogenic factors, the predictive uncertainty

501   limits are also uncertain, which means that we are unable to estimate predictive uncertainty

502   limits even if our estimation of modeling uncertainty limits are accurate. Hence in application

503   of uncertainty analysis for hydrologic prediction, how to extend modeling uncertainty limits

504   to predictive uncertainty limits remains a challenge for applying hydrologic models to water

505   resources-related management and design problems.

506   **4. CONCLUSIONS**

507       In this paper, we presented the application of GA and BMA to simultaneously conduct

508   calibration and uncertainty analysis of SWAT. The methodology provides a practical and

509   flexible tool for jointly using multiple model structures within the SWAT model system. This

510   method was tested in two basins. In the LREB, we selected six SWAT models with different

511   evapotranspiration and flow routing algorithms, and tested this method using daily

512   streamflow. In the YRHB, we selected three SWAT models with different snow routing

513   modules, and used monthly streamflow data to test this method. The test results show that

514   this combined method can provide deterministic predictions better than or comparable to the

21

515 best calibrated model using GA. Further inspection of the 66.7% and 90% uncertainty

516 intervals show that the combination of GA and BMA can provide reasonable uncertainty

517 estimation. The differences between the computed percentage of coverage values and the

518 corresponding expected coverage percentages are within 10% for both calibration and

519 validation periods in these two test basins. It is anticipated that the combination of GA and

520 BMA methods will have significant implications related to policy development. The method

521 reduces the uncertainty associated with selecting any single model, thereby increasing the

522 level of confidence in the simulation results. This is a critical component of policy

523 assessments which are based upon modeling results and one which will become more routine

524 in the future.

525

526

## Acknowledgements

527

531 **Appendix A:**

532 1. Initialization:

533 Set $Iter = 0$, $w_k^{Iter} = \dfrac{1}{K}$, and $\sigma^2_{k,Iter} = \dfrac{1}{K}\sum_{t=1}^{T}(\sum_{k=1}^{K}(y_t - f_{k,t})^2 / T)$, and fit the regression

534 coefficients $a_k$ and $b_k$ for each candidate model using linear regression.

535 where $T$ is the total number of data points in the calibration period, and $Iter$ is the iteration
536 number.

537 2. Computing the initial likelihood:

538
$$\ell(\mathbf{\theta}^{Iter}) = \log\left(\sum_{k=1}^{K} w_k g(y \mid f_k)\right) \qquad\qquad A1$$

539 where $g(y \mid f_k)$ is calculate as $\sum_{t=1}^{T} g(y_t \mid f_{k,t}, \sigma^2_{k,Iter})$. $g(y_t \mid f_{k,t}, \sigma^2_{k,Iter})$ represents a normal

540 distribution center at $a_k + b_k f_{k,t}$ with variance of $\sigma^2_{k,Iter}$

541 3. Executing the expectation step

542 Set $Iter = Iter + 1$

543 For $k = 1,2,\ldots,K$ and $t = 1,2,\ldots,T$, $\hat{z}^{Iter}_{k,t} = g(y_t \mid f_{k,t}, \sigma_{k,Iter-1}) \Big/ \sum_{k=1}^{K} g(y_t \mid f_{k,t}, \sigma_{k,Iter-1})$

544 4. Executing the maximization step

545 Compute the weight for each model: $w_k^{Iter} = \dfrac{1}{T}\sum_{t=1}^{T}\hat{z}^{Iter}_{k,t}$

546 Update the variance of each model: $\sigma^2_{k,Iter} = \sum_{t=1}^{T}\hat{z}^{Iter}_{k,t}(y_t - f_{k,t})^2 \Big/ \sum_{t=1}^{T}\hat{z}^{Iter}_{k,t}$

547 5. Update the likelihood $\ell(\mathbf{\theta}^{Iter})$ using equation A1

548 6. Checking convergence:

549 If $\ell(\mathbf{\theta}^{Iter}) - \ell(\mathbf{\theta}^{Iter-1})$ is less than or equal to a pre-specified tolerance level ($10^{-6}$), stop; else go
550 back to Step 3.

551  **References**

552  Abbaspour, K.C., Johnson, C. A., van Genuchten, M.T., 2004. Estimating uncertain flow and
553      transport parameters using a sequential uncertainty fitting procedure. *Vadose Zone*
554      *Journal* 3(4), 1340–1352.

555  Ajami, K., Duan, Q., Gao, X., Sorooshian, S., 2006. Multi-model combination techniques for
556      hydrological forecasting: application to distributed model intercomparison project results.
557      *Journal of Hydrometeorology* 8, 755–768.

558  Arabi, M., Govindaraju, R.S., Hantush M. M., 2007. A probabilistic approach for analysis of
559      uncertainty in the evaluation of watershed management practices. *Journal of Hydrology*
560      333, 459–471.

561  Arnold, J.G., Srinivasan, R., Muttiah, R.S., Williams, J.R., 1998. Large-area hydrologic
562      modeling and assessment: Part I. Model development. *Journal of the American Water*
563      *Resources Association* 34(1): 73-89.

564  Bekele, G. E., Nicklow,W.J., 2007. Multi-objective automatic calibration of SWAT using
565      NSGA-II. *Journal of Hydrology* 341: 165-176.

566  Beven, K.J. 2006. A manifesto for the enquiringly thesis, *Journal of Hydrology* 320, 18-36.

567  Beven, K., Binley, A., 1992. The future of distributed models – model calibration and
568      uncertainty prediction. *Hydrological Processes* 6 (3), 279–298.

569  Beven, K.J., 2000. *Rainfall-runoff Modeling: The Primer*. John Wiley & Sons Press: New
570      York.

571  Beven, K.J., Freer, J., 2001. Equifinality, data assimilation, and uncertainty estimation in
572      mechanistic modeling of complex environmental systems. *Journal of Hydrology* 249, 11-
573      29.

574  Bosch, D.D., Sheridan, J.M., Lowrance, R.R., Hubbard, R.K., Strickland, T.C., Feyereisen,
575      G.W., Sullivan, D.G., 2007. Little River Experimental Watershed database. *Water*
576      *Resources Research 43*, W09470, doi:10.1029/2006WR005844.

577  Bosch, D.D., Sullivan, D.G., Sheridan, J. M., 2006. Hydrologic impacts of land-use changes
578      in coastal plain watersheds. *Transactions of the ASABE* 49(2): 423−432.

579  Butts, M.B., Payne, J.T., Kristensen, M. and Madsen, H., 2004a. An evaluation of the impact
580      of model structure on hydrological modelling uncertainty for streamflow simulation.
581      *Journal of Hydrology* 298: 222-241.

582  Butts, M.B., Payne, J.T. and Overgaard, J., 2004b. Improving streamflow simulations and
583      flood forecasting with multimodel ensemble. In: P.B. Liong (Editor), 6th International
584      Conference on Hydroinformatics. World Scientific Publishing, Singapore.

585  Chakravarti I.M., Laha R.G., and Roy J. 1967. *Handbook of Methods of Applied Statistics*.
586      John Wiley and Sons, New York, USA.

587  Duan, Q., Ajami, N. K., Gao, X., Sorooshian, S., 2007. Multi-model ensemble hydrologic
588      prediction using Bayesian model averaging. *Advances in Water Resources* 30(5), 1371-
589      1386.

590   Duan, Q., Sorooshian, S., Gupta, V.K., 1992. Effective and efficient global optimization for
591       conceptual rainfall-runoff models. *Water Resources Research* 28(4): 1015-1031.

592   Fontaine, T.A., Cruickshank, T.S., Arnold, J.G., Hotchkiss, R.H., 2002. Development of a
593       snowfall-snowmelt routine for mountainous terrain for the soil water assessment tool
594       (SWAT). *Journal of Hydrology* 262, 209-223.

595   Gassman, P.W., Reyes, M., Green, C.H., Arnold, J.G., 2007. The Soil and Water Assessment
596       Tool: Historical development, applications, and future directions. *Transactions of the*
597       *ASABE* 50(4): 1212-1250.

598   Gelman, A., Carlin J. B., Stern H. S., Rubin D. B., 2003. Bayesian Data Analysis (2nd
599       edition). Chapman & Hall/CRC: Boca Raton, Florida, USA.

600   Gelman, A. Rubin D. B., 1995. Avoiding model selection in Bayesian social research.
601       Sociological Methodology 25: 165-173.

602   Goldberg, D., 1989. *Genetic Algorithms in Search, Optimization and Machine Learning*.
603       Addison-Wesley, Reading, Massachusetts, USA.

604   Gupta, H. V., Sorooshian, S., Yapo, P.O., 1998. Toward improved calibration of hydrologic
605       models: Multiple and noncommensurate measures of information. *Water Resources*
606       *Research* 34(4): 751–763.

607   Gupta, H.V., Sorooshian, S., Yapo, P. O., 1999. Status of automatic calibration for
608       hydrologic models: Comparison with multilevel expert calibration. *Journal of Hydrologic*
609       *Engineering* 4(2), 135-143.

610   Harmel, R.D., Smith, P. K., 2007. Consideration of measurement uncertainty in the
611       evaluation of goodness-of-fit in hydrologic and water quality modeling. *Journal of*
612       *Hydrology* 337, 326–336.

613   Holland, J., 1975. *Adaptation in Natural and Artificial Systems*. University of Michigan
614       Press, Ann Arbor, Michigan, USA.

615   Hsu, K., Moradkhani, H., Sorooshian, S., 2009. A sequential Bayesian approach for
616       hydrologic model selection and prediction. *Water Resources Research*
617       doi:10.1029/2008WR006824.

618   Kannan, N., Santhi, C., Arnold, J.G., 2008. Development of an automated procedure for
619       estimation of the spatial variation of runoff in large river basins. *Journal of Hydrology*
620       359, 1–15.

621   Kavetski, D., Kuczera G., Franks S. W., 2006. Bayesian analysis of input uncertainty in
622       hydrological modeling: 1. Theory. *Water Resources Research 42*, W03407,
623       doi:10.1029/2005WR004368.

624   Kuczera, G., Kavetski, D., Franks, S., Thyer, M. 2006. Towards a Bayesian total error
625       analysis of conceptual rainfall-runoff models: Characterising model error using storm-
626       dependent parameters. *Journal of Hydrology* 331, 161– 177.

627   Kuczera, G., Parent, E., 1998. Monte Carlo assessment of parameter uncertainty in
628       conceptual catchment models: the Metropolis algorithm. *Journal of Hydrology* 211, 69-
629       85.

630 Legates, D.R., McCabe, G.J., 1999. Evaluating the use of "goodness of fit" measures in
631   hydrologic and hydroclimatic model validation. *Water Resources Research* 35(1): 233-
632   241.

633 Liu, Y., Gupta,V., 2007. Uncertainty in hydrologic modeling: Toward an integrated data
634   assimilation framework. *Water Resources Research* 43, W07401,
635   doi:10.1029/2006WR005756.

636 Montgomery, J., Nyhan, B., 2008. *Bayesian Model Averaging: Theoretical developments*
637   *and practical applications*. Available at http://www.duke.edu/~bjn3/montgomery-nyhan-
638   bma.pdf. Accessed on Oct. 8, 2008.

639 Moriasi, D.N., Arnold, J.G., Van Liew, M.W., Bingner, R.L., Harmel, R. D., Veith, T.L.,
640   2007. Model evaluation guidelines for systematic quantification of accuracy in watershed
641   simulations. *Transactions of the ASABE* 50(3), 885−900.

642 Neitsch, S.L., Arnold, J.G., Kiniry, J.R., King, K.W., Williams, J.R., 2005a. Soil and Water
643   Assessment Tool (SWAT) theoretical documentation. Blackland Research Center, Texas
644   Agricultural Experiment Station, Temple, Texas, BRC Report 02-05.

645 Neitsch, S.L., Arnold, J.G., Kiniry, J.R., Srinivasan, R., Williams, J. R., 2005b. *Soil and*
646   *Water Assessment Tool (SWAT) users manual*. Blackland Research Center, Texas
647   Agricultural Experiment Station, Temple, Texas, BRC Report 02-06.

648 Neuman, S. P., 2003. Maximum Likelihood Bayesian averaging of uncertain model
649   predictions. *Stochastic Environmental Research and Risk Assessment* 17: 291-305.

650 Raftery, A.E., Gneiting, T., Balabdaoui, F., Polakowski, M., 2005. Using bayesian model
651   averaging to calibrate forecast ensembles. *Monthly Weather Review* 113: 1155–1174.

652 Raftery, A.E., Zheng, Y., 2003. Discussion: performance of Bayesian model averaging.
653   *Journal of the American Statistical Association* 98 (464), 931–938.

654 Reca, J., Martínez, J., 2006. Genetic algorithms for the design of looped irrigation water
655   distribution networks. *Water Resources Research* 42, W05416,
656   doi:10.1029/2005WR004383.

657 Santhi C., Arnold, J.G., Williams, J.R., Dugas, W.A., Hauck L., 2001. Validation of the
658   SWAT model on a large river basin with point and nonpoint sources. *Journal of the*
659   *American Water Resources Association* 37(5), 1169-1188.

660 Schaffer, J.D., Caruana, R.A., Eshelman, L.J., Das, R., 1989. A study of control parameters
661   affecting online performance of genetic algorithms for function optimization. In:
662   *Proceedings of the Third International Conference on Genetic algorithms* (ed. By
663   Schaffer, J. D.), 51-60. Morgan Kaufmann, San Mateo, California, USA.

664 Sheridan, J.M., 1997. Rainfall-streamflow relations for coastal plain watersheds.
665   *Transactions of ASAE* 13(3): 333-344.

666 Shirmohammadi, A., Chaubey, I., Harmel, R.D. Bosch, D.D. Munoz-Carpena, R.C.
667   Dharmasi, A. Arabi, S.M., Wolfe, M.L. Frankenberger, J., Graff, C., Sohrabi. T.M., 2006.
668   Uncertainty in TMDL Models. *Transactions of the ASABE* 49(4):1033-1049.

669  Van Griensven, A., Meixner, T., 2006. Methods to quantify and identify the sources of
670      uncertainty for river basin water quality models. *Water Science Technology* 53 (1), 51–
671      59.

672  Van Griensven, A., Meixner, T., Grunwald, S., Bishop, T., Di luzio, M., Srinivasan, R., 2006.
673      A global sensitivity analysis tool for the parameters of multi-variable catchment models.
674      *Journal Hydrology* 324: 10-23.

675  Van Griensven, A., Meixner, T., Srinivasan, R., Grunwals, S., 2008. Fit-for-purpose analysis
676      of uncertainty using split-sampling evaluations. *Hydrological Sciences Journal* 53(5),
677      1090-1103.

678  Van Liew, M.W., Arnold, J.G., Bosch, D.D., 2005. Problems and Potential of Autocalibrating
679      a Hydrologic Model. *Transactions of the ASAE* 48(3), 1025-1040

680  Van Liew, M. W., Veith, T. L., Bosch, D. D., Arnold, J. G., 2007. Suitability of SWAT for
681      the Conservation Effects Assessment Project: A comparison on USDA ARS watersheds.
682      *Journal of Hydrologic Engineering* 12(2), 173-189.

683  Vrugt, J.A., Gupta, H.V., Bouten, W., Sorooshian, S., 2003. A shuffled complex evolution
684      metropolis algorithm for optimization and uncertainty assessment of hydrologic model
685      parameters. *Water Resources Research* 39(8):1201, doi:10.1029/2002WR001642.

686  Vrugt, J. A., Nuallain B., Robinson B. A., Bouten W., Dekker S.C., Sloot P. M. A., 2006.
687      Application of parallel computing to stochastic parameter estimation in environmental
688      models. Computers & Geosciences, 32(8), 1139 - 1155

689  Vrugt, J. A., Robinson, B. A., 2007. Treatment of uncertainty using ensemble methods:
690      Comparison of sequential data assimilation and Bayesian model averaging. *Water
691      Resources Research* 43, W01411, doi:10.1029/2005WR004838.

692  Yang, J., Reichert, P., Abbaspour, K. C., Xia, J., Yang, H., 2008. Comparing uncertainty
693      analysis techniques for a SWAT application to the Chaohe Basin in China. *Journal of
694      Hydrology* 358, 1–23.

695  Yang, K., Reichert, P., Abbaspour, K. C., Yang, H., 2007. Hydrological modelling of the
696      Chaohe Basin in China: Statistical model formulation and Bayesian inference. *Journal of
697      Hydrology* 340, 167–182.

698  Zhang, X, Srinivasan, R, Debele, B., Hao, F., 2008a. Runoff simulation of the Headwaters of
699      the Yellow River using the SWAT model with three snowmelt algorithms. *Journal of the
700      American Water Resources Association* 44(1), 48-61. DOI: 10.1111 ∕ j.1752-
701      1688.2007.00137.x.

702  Zhang, X., 2008b. Evaluating and developing parameter optimization and uncertainty
703      analysis methods for a computationally intensive distributed hydrologic model. Ph. diss.
704      Texas A&M University, College Station, Texas, USA.

705  Zhang, X., Liang, F., Srinivasan, R., Van Liew, M., 2009a. Estimating Uncertainty of
706      Streamflow Simulation using Bayesian Neural Networks. *Water Resources Research*
707      doi:10.1029/2008WR007030.

708 Zhang, X., Srinivasan, R., Van Liew, M., 2009b. Approximating the SWAT Model Using
709     Artificial Neural Network and Support Vector Machine. *Journal of the American Water*
710     *Resources Association* 45(2): 460-474.

711 Zhang, X., Srinivasan, R., Zhao, K., Van Liew, M., 2009c. Evaluation of global optimization
712     algorithms for parameter calibration of a computationally intensive hydrologic model.
713     *Hydrological Processes* 23(3): 430-441.

714 **List of Tables**

715

725

726

727

728                          Table 1 Parameters for calibration in SWAT model.

|    | Parameter | Description | Range |
|----|-----------|-------------|-------|
| 1  | CN2       | Curve Number | ±20% |
| 2  | ESCO      | Soil Evaporation compensation factor | 0-1 |
| 3  | SOL_AWC   | Available soil water capacity | ±20% |
| 4  | GW_REVAP  | Ground water re-evaporation coefficient | 0.02-0.2 |
| 5  | REVAPMN   | Threshold depth of water in the shallow aquifer for re-evaporation to occur (mm). | 0-500 |
| 6  | GWQMN     | Threshold depth of water in the shallow aquifer required for return flow to occur (mm) | 0-5000 |
| 7  | GW_DELAY  | Groundwater delay (days) | 0-50 |
| 8  | ALPHA_BF  | Base flow recession constant | 0-1 |
| 9  | RCHRG_DP  | Deep aquifer percolation fraction | 0-1 |
| 10 | CH_K2     | Effective hydraulic conductivity in main channel alluvium (mm/hr) | 0.01-150 |
| 11 | SURLAG    | Surface runoff lag coefficient (day) | 0-10 |

729      Table 2 Evaluation coefficients for the six SWAT models, arithmetic mean, and BMA mean
730          in the LREB for both calibration and validation periods.

| Coefficients / Models | Calibration | | | Validation | | |
|---|---|---|---|---|---|---|
| | PBIAS | $R^2$ | NSE | PBIAS | $R^2$ | NSE |
| SWAT-HG-MK | -0.72% | 0.76 | 0.74 | -8.24% | 0.82 | 0.71 |
| SWAT-HG-VS | 6.66% | 0.76 | 0.75 | 27.07% | 0.81 | 0.76 |
| SWAT-PM-MK | 23.90% | 0.77 | 0.76 | 35.13% | 0.82 | 0.8 |
| SWAT-PM-VS | 24.04% | 0.72 | 0.71 | 39.77% | 0.8 | 0.75 |
| SWAT-PT-MK | 11.26% | 0.79 | 0.78 | 23.49% | 0.85 | 0.74 |
| SWAT-PT-VS | 22.94% | 0.71 | 0.7 | 46.82% | 0.77 | 0.5 |
| Ensemble Mean | 14.60% | **0.81** | 0.79 | 27.34% | 0.86 | 0.84 |
| BMA mean | **0.00%** | **0.81** | **0.81** | **3.07%** | **0.87** | **0.86** |

731    Table 3 Calibrated parameter values for the six models in LREB.

| Model / Parameter | SWAT_HG_MK | SWAT_HG_VS | SWAT_PM_MK | SWAT_PM_VS | SWAT_PT_MK | SWAT_PT_VS |
|---|---|---|---|---|---|---|
| CN | 9% | -17% | 8% | -17% | 6% | 20% |
| ESCO | 0.46 | 0.89 | 0.88 | 0.91 | 0.38 | 0.78 |
| Surlag | 9.99 | 2.8 | 9.78 | 1.1 | 9.69 | 2.3 |
| ALPHA_BF | 0.23 | 0.61 | 0.17 | 0.45 | 0.37 | 0.55 |
| GW_REVAP | 0.15 | 0.15 | 0.2 | 0.2 | 0.08 | 0.1 |
| SOL_AWC | 7% | -20% | 18% | 16% | 18% | 25% |
| CH_K2 | 144 | 147 | 146 | 130 | 131 | 147 |
| GW_DELAY | 22.57 | 3.7 | 18.87 | 2.19 | 22.8 | 3.07 |
| RCHRG_DP | 0.79 | 0.01 | 0.66 | 0.45 | 0.33 | 0.68 |
| GWQMN | 9.16 | 103.44 | 45.91 | 103.69 | 95.14 | 168.81 |
| REVAPMN | 215.14 | 24.59 | 486.46 | 402.32 | 263.62 | 190.1 |

732

733

734    Table 4 Evaluation coefficients for the three SWAT models, arithmetic mean, and BMA
735    mean in the YRHB for both calibration and validation periods.

| Coefficients / Models | Calibration | | | Validation | | |
|---|---|---|---|---|---|---|
| | PBIAS | $R^2$ | NSE | PBIAS | $R^2$ | NSE |
| SWAT-DD | -17.71% | 0.82 | 0.77 | -17.98% | 0.84 | 0.78 |
| SWAT-ELEV | -4.63% | 0.85 | 0.84 | **-0.31%** | 0.83 | 0.83 |
| SWAT-SNOW17 | 4.76% | 0.87 | 0.84 | 7.12% | 0.85 | 0.78 |
| Ensemble Mean | -5.86% | **0.88** | 0.87 | -3.72% | **0.87** | **0.87** |
| BMA Mean | **0.00%** | **0.88** | **0.88** | 3.71% | **0.87** | **0.87** |

736

737

738

32

739 Table 5 Calibrated parameter values for the three models in YRHB.

| Parameter / Model | CN | ESCO | Surlag | ALPHA_BF | GW_REVAP |
|---|---|---|---|---|---|
| SWAT-DD | 14% | 0.28 | 4.90 | 0.16 | 0.03 |
| SWAT-ELEV | 7% | 0.36 | 3.80 | 0.33 | 0.04 |
| SWAT_SNOW17 | 2% | 0.18 | 7.40 | 0.51 | 0.08 |

740

741

742 Table 6 Evaluation coefficients obtained using different number of candidate models in BMA
743 in the LREB.

| Coefficients / Number of Candidate models | Calibration | | | | | Validation | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | PBIAS | $R^2$ | NSE | 66.7% POC | 90% POC | PBIAS | $R^2$ | NSE | 66.7% POC | 90% POC |
| 6 | 0.00% | 0.81 | 0.81 | 76.04% | 91.14% | 3.07% | 0.87 | 0.86 | 74.41% | 96.53% |
| 5 | 0.00% | 0.81 | 0.81 | 75.34% | 91.32% | 1.02% | 0.86 | 0.86 | 74.15% | 92.93% |
| 4 | 0.00% | 0.81 | 0.8 | 74.33% | 90.99% | 1.32% | 0.86 | 0.86 | 73.89% | 94.94% |
| 3 | 0.00% | 0.80 | 0.8 | 73.96% | 91.71% | 3.21% | 0.86 | 0.85 | 73.36% | 94.01% |
| 2 | 0.00% | 0.80 | 0.8 | 75.89% | 93.15% | 10.58% | 0.85 | 0.84 | 77.02% | 95.08% |

744

745

746

33

## List of Figures

**Figure 1**

**Figure 4**

**Figure 5**

**Figure 8**



Figure 8. Streamflow (cms) — Monthly flow from January 1976 to December 1985. Top panel: 66.7% interval with observed data. Bottom panel: 90% interval with observed data.

**Figure 9**



Monthly flow from January 1986 to December 1990