



Pianosi, F. (2015). Curses, Tradeoffs, and Scalable Management: Advancing Evolutionary Multiobjective Direct Policy Search to Improve Water Reservoir Operations. ASCE Journal of Water Resources Planning and Management, 142(2), [04015050]. 10.1061/(ASCE)WR.1943-5452.0000570

Peer reviewed version

Link to published version (if available): 10.1061/(ASCE)WR.1943-5452.0000570

Link to publication record in Explore Bristol Research PDF-document

# **University of Bristol - Explore Bristol Research**

#### **General rights**

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available: http://www.bristol.ac.uk/pure/about/ebr-terms.html

### Take down policy

Explore Bristol Research is a digital archive and the intention is that deposited content should not be removed. However, if you believe that this version of the work breaches copyright law please contact open-access@bristol.ac.uk and include the following information in your message:

- Your contact details
- Bibliographic details for the item, including a URL
- An outline of the nature of the complaint

On receipt of your message the Open Access Team will immediately investigate your claim, make an initial judgement of the validity of the claim and, where appropriate, withdraw the item in question from public view.

# Curses, tradeoffs, and scalable management: advancing evolutionary multi-objective direct policy search to improve water reservoir operations

<sup>4</sup> Matteo Giuliani<sup>1</sup>; Andrea Castelletti<sup>23</sup>; Francesca Pianosi<sup>4</sup>; Emanuele Mason<sup>5</sup>; Patrick M. Reed<sup>6</sup>

## 5 ABSTRACT

Optimal management policies for water reservoir operation are generally designed via 6 stochastic dynamic programming (SDP). Yet, the adoption of SDP in complex real-world 7 problems is challenged by the three curses of dimensionality, of modeling, and of multiple 8 objectives. These three curses considerably limit SDP's practical application. Alterna-9 tively, in this study, we focus on the use of evolutionary multi-objective direct policy search 10 (EMODPS), a simulation-based optimization approach that combines direct policy search. 11 nonlinear approximating networks and multi-objective evolutionary algorithms to design 12 Pareto approximate closed-loop operating policies for multi-purpose water reservoirs. Our 13 analysis explores the technical and practical implications of using EMODPS through a care-14 ful diagnostic assessment of the effectiveness and reliability of the overall EMODPS solution 15 design as well as of the resulting Pareto approximate operating policies. The EMODPS 16 approach is evaluated using the multi-purpose Hoa Binh water reservoir in Vietnam, where 17 water operators are seeking to balance the conflicting objectives of maximizing hydropower 18

<sup>&</sup>lt;sup>1</sup>PhD, Dept. of Electronics, Information, and Bioengineering, Politecnico di Milano, P.za Leonardo da Vinci, 32, 20133 Milano, Italy. E-mail: matteo.giuliani@polimi.it.

<sup>&</sup>lt;sup>2</sup>Associate Professor, M. ASCE, Dept. of Electronics, Information, and Bioengineering, Politecnico di Milano, P.za Leonardo da Vinci, 32, 20133 Milano, Italy. E-mail: andrea.castelletti@polimi.it.

<sup>&</sup>lt;sup>3</sup>Senior scientist, Institute of Environmental Engineering, ETH Zurich, Ramistrasse 101, 8092 Zurich Switzerland.

<sup>&</sup>lt;sup>4</sup>Research Associate, Dept. of Civil Engineering, University of Bristol, Queen's Building, University Walk, Bristol BS8 1TR, UK. Email: francesca.pianosi@bristol.ac.uk.

<sup>&</sup>lt;sup>5</sup>PhD student, Dept. of Electronics, Information, and Bioengineering, Politecnico di Milano, P.za Leonardo da Vinci, 32, 20133 Milano, Italy. E-mail: emanuele.mason@polimi.it.

<sup>&</sup>lt;sup>6</sup>Professor, M. ASCE, School of Civil and Environmental Engineering, University of Cornell, 211 Hollister Hall, Ithaca, USA.E-mail: patrick.reed@cornell.edu.

production and minimizing flood risks. A key choice in the EMODPS approach is the selec-19 tion of alternative formulations for flexibly representing reservoir operating policies. In this 20 study, we distinguish the relative performance of two widely used nonlinear approximating 21 networks, namely Artificial Neural Networks and Radial Basis Functions. Our results show 22 that RBF solutions are more effective than ANN ones in designing Pareto approximate poli-23 cies for the Hoa Binh reservoir. Given the approximate nature of EMODPS, our diagnostic 24 benchmarking uses SDP to evaluate the overall quality of the attained Pareto approximate 25 results. Although the Hoa Binh test case's relative simplicity should maximize the potential 26 value of SDP, our results demonstrate that EMODPS successfully dominates the solutions 27 derived via SDP. 28

<sup>29</sup> Keywords: water management, direct policy search, multi-objective evolutionary algorithm

#### **30 INTRODUCTION**

Climate change and growing populations are straining freshwater availability worldwide 31 (McDonald et al. 2011) to the point that many large storage projects are failing to produce 32 the level of benefits that provided the economic justification for their development (Ansar 33 et al. 2014). In a rapidly changing context, operating existing infrastructures more effi-34 ciently, rather than planning new ones, is a critical challenge to balance competing demands 35 and performance uncertainties (Gleick and Palaniappan 2010). Yet, most major reservoirs 36 have had their operations defined in prior decades (U.S. Army Corps of Engineers 1977; 37 Loucks and Sigvaldason 1982), assuming "normal" hydroclimatic conditions and consider-38 ing a restricted number of operating objectives. The effectiveness of these rules is however 39 limited, as they are not able to adapt release decisions when either the hydrologic system 40 deviates from the assumed baseline conditions or additional objectives emerge over time. 41 On the contrary, closing the loop between operational decisions and evolving system condi-42 tions provides the adaptive capacity needed to face growing water demands and increasingly 43 uncertain hydrologic regimes (Soncini-Sessa et al. 2007). 44

<sup>45</sup> In the literature, the design problem of closed-loop operating policies for managing water

storages has been extensively studied since the seminal work by Rippl (1883). From the first 46 applications by Hall and Buras (1961), Maass et al. (1962), and Esogbue (1989), Dynamic 47 Programming (DP) and its stochastic extension (SDP) are probably the most widely used 48 methods for designing optimal operating policies for water reservoirs (for a review, see Yeh 49 (1985); Labadie (2004); Castelletti et al. (2008), and references therein). SDP formulates the 50 operating policy design problem as a sequential decision-making process, where a decision 51 taken now produces not only an immediate reward, but also affects the next system state 52 and, through that, all the subsequent rewards. The search for optimal policies relies on the 53 use of value functions defined over a discrete (or discretized) state-decision space, which are 54 obtained by looking ahead to future events and computing a backed-up value. In principle, 55 SDP can be applied under relatively mild modeling assumptions (e.g., finite domains of state, 56 decision and disturbance variables, time-separability of objective functions and constraints). 57 In practice, the adoption of SDP in complex real-world water resources problems is challenged 58 by three curses that considerably limit its use, namely the curse of dimensionality, the curse 59 of modeling, and the curse of multiple objectives. 60

The curse of dimensionality, first introduced by Bellman (1957), means that the compu-61 tational cost of SDP grows exponentially with the state vector dimensionality. SDP would 62 be therefore inapplicable when the dimensionality of the system exceeds 2 or 3 storages 63 (Loucks et al. 2005). In addition, particularly in such large systems, the disturbances (e.g., 64 inflows) are likely to be both spatially and temporally correlated. While including space 65 variability in the identification of the disturbance's probability distribution function (pdf) 66 can be sometimes rather complicated, it does not add to SDP's computational complex-67 ity. Alternatively, properly accounting for temporal correlation requires using a dynamic 68 stochastic model, which contributes additional state variables and exacerbates the curse of 69 dimensionality. 70

The curse of modeling was defined by Tsitsiklis and Van Roy (1996) to describe the SDP requirement that, in order to solve the sequential decision-making process at each

stage in a step-based optimization, any information included into the SDP framework must 73 be explicitly modeled to fully predict the one-step ahead model transition used for the 74 estimation of the value function. This information can be described either as a state variable 75 of a dynamic model or as a stochastic disturbance, independent in time, with an associated 76 pdf. As a consequence, exogenous information (i.e., variables that are observed but are 77 not affected by the decisions, such as observations of inflows, precipitation, snow water 78 equivalent, etc.), which could potentially improve the reservoir operation (Tejada-Guibert 79 et al. 1995; Faber and Stedinger 2001), cannot be explicitly considered in conditioning the 80 decisions, unless a dynamic model is identified for each additional variable, thus adding to the 81 curse of dimensionality (i.e., additional state variables). Moreover, SDP cannot be combined 82 with high-fidelity process-based simulation models (e.g., hydrodynamic and ecologic), which 83 require a warm-up period and cannot be employed in a step-based optimization mode. 84

The curse of multiple objectives (Powell 2007) is related to the generation of the full set 85 of Pareto optimal (or approximate) solutions to support a posteriori decision making (Cohon 86 and Marks 1975) by exploring the key alternatives that compose system tradeoffs, providing 87 decision makers with a broader context where their preferences can evolve and be exploited 88 opportunistically (Brill. et al. 1990; Woodruff et al. 2013). Most of the DP-family methods 89 relies on single-objective optimization algorithms, which require a scalarization function (e.g., 90 convex combination or non-linear Chebyshev scalarization) to reduce the dimensionality of 91 the objective space to a single-objective problem (Chankong and Haimes 1983; ReVelle and 92 McGarity 1997). The single-objective optimization is then repeated for every Pareto optimal 93 point generated by using different scalarization values (Soncini-Sessa et al. 2007). However, 94 this process is computationally very demanding in many-objective optimization problems, 95 namely when the number of objectives grows to three or more (Fleming et al. 2005), and the 96 accuracy in the approximation of the Pareto front might be degraded given the non-linear 97 relationships between the scalarization values and the corresponding objectives values. 98

Approximate Dynamic Programming (Powell 2007) and Reinforcement Learning (Buso-

niu et al. 2010) seek to overcome some or all the SDP curses through three different ap-100 proaches: (i) value function-based methods, which compute an approximation of the value 101 function (Bertsekas and Tsitsiklis 1996); (ii) on-line methods, which rely on the sequential 102 resolution of multiple open-loop problems defined over a finite receding horizon (Bertsekas 103 2005); (iii) policy search-based methods, which use a simulation-based optimization to itera-104 tively improve the operating policies based on the simulation outcome (Marbach and Tsitsik-105 lis 2001). However, the first two approaches still rely on the estimation (or approximation) 106 of the value function with single-objective optimization algorithms. Simulation-based op-107 timization, instead, represents a promising alternative to reduce the limiting effects of the 108 three curses of SDP by first parameterizing the operating policy using a given family of 100 functions and, then, by optimizing the policy parameters (i.e., the decision variables of the 110 problem) with respect to the operating objectives of the problem. This approach is generally 111 named direct policy search (DPS, see Rosenstein and Barto (2001)) and is also known in 112 the water resources literature as parameterization-simulation-optimization by Koutsoyiannis 113 and Economou (2003), where has been adopted in several applications (Guariso et al. 1986; 114 Oliveira and Loucks 1997; Cui and Kuczera 2005; Dariane and Momtahen 2009; Guo et al. 115 2013). 116

The simulation-based nature of DPS offers some advantages over the DP-family methods. 117 First, the variable domain does not need to be discretized, thus reducing the curse of dimen-118 sionality. The complexity of the operating policy (i.e., the number of policy inputs/outputs) 119 however depends on the dimensionality of the system. The higher the number of reservoirs, 120 the more complex is the policy to design, which requires a large number of parameters. 121 Second, DPS can be combined with any simulation model and does not add any constraint 122 on modeled information, allowing the use of exogenous information in conditioning the de-123 cisions. Third, when DPS problems involve multiple objectives, they can be coupled with 124 truly multi-objective optimization methods, such as multi-objective evolutionary algorithms 125 (MOEAs), which allow estimating an approximation of the Pareto front in a single run of 126

127 the algorithm.

Following Nalbantis and Koutsoyiannis (1997), DPS can be seen as an optimization-based 128 generalization of well known simulation-based, single-purpose heuristic operating rules (U.S. 129 Army Corps of Engineers 1977). The New York City rule (Clark 1950), the spill-minimizing 130 "space rule" (Clark 1956), or the Standard Operating Policy (Draper and Lund 2004) can 131 all be seen as parameterized single-purpose policies. Many of these rules are based largely 132 on empirical or experimental successes and they were designed, mostly via simulation, for 133 single-purpose reservoirs (Lund and Guzman 1999). In more complex systems, such as 134 networks of multi-purpose water reservoirs, the application of DPS is more challenging due 135 to the difficulties of choosing an appropriate family of functions to represent the operating 136 policy. Since DPS can, at most, find the best possible solution within the prescribed family 137 of functions, a bad approximating function choice can strongly degrade the final result. For 138 example, piecewise linear approximations have been demonstrated to work well for specific 139 problems, such as hedging rules or water supply (Oliveira and Loucks 1997). In other 140 problems (e.g., hydropower production), the limited flexibility of these functions can however 141 restrict the search to a subspace of policies that, likely, does not contain the optimal one. In 142 many cases, the choice of the policy architecture can not be easily inferred either from the 143 experience of the water managers, who may not be operating the system at full attainable 144 efficiency, or a priori on the basis of empirical considerations, when the system is under 145 construction and data about the historical regulation are not yet available. A more flexible 146 function, depending on a larger number of parameters, has hence to be selected to ensure 147 the possibility of approximating the unknown optimal solution of the problem to any desired 148 degree of accuracy. In this work, we have adopted two widely used nonlinear approximating 149 networks (Zoppoli et al. 2002), namely Artificial Neural Networks (ANNs) and Radial Basis 150 Functions (RBFs), which have been demonstrated to be universal approximators under mild 151 assumptions on the activation functions used in the hidden layer (for a review see Tikk et al. 152 (2003) and references therein). 153

The selected policy parameterization strongly influences the selection of the optimiza-154 tion approach, which is often case study dependent and may require ad-hoc tuning of the 155 optimization algorithms. Simple parameterizations, defined by a limited number of param-156 eters, can be efficiently optimized via ad-hoc gradient-based methods (Peters and Schaal 157 2008; Sehnke et al. 2010). On the contrary, gradient-free global optimization methods are 158 preferred when the complexity of the policy parameterization, and consequently the number 159 of parameters to optimize, increases. In particular, evolutionary algorithms (EAs) have been 160 successfully applied in several policy search problems characterized by high-dimensional de-161 cision spaces as well as noisy and multi-modal objective functions (Whitley et al. 1994; 162 Moriarty et al. 1999; Whiteson and Stone 2006; Busoniu et al. 2011). Indeed, EAs search 163 strategies, which are based on ranking of candidate solutions, better handle the performance 164 uncertainties than methods relying on the estimation of absolute performance or perfor-165 mance gradient (Heidrich-Meisner and Igel 2008). This property is particular relevant given 166 the stochasticity of water resources systems. In this work, we address the challenges posed 167 by multi-objective optimization under uncertainty by using the self-adaptive Borg MOEA 168 (Hadka and Reed 2013). The Borg MOEA has been shown to be highly robust across a 169 diverse suite of challenging multi-objective problems, where it met or exceeded the perfor-170 mance of other state-of-the-art MOEAs (Reed et al. 2013). In particular, the Borg MOEA 171 overcomes the limitations of tuning the algorithm parameters to the specific problems by 172 employing multiple search operators, which are adaptively selected during the optimization 173 based on their demonstrated probability of generating quality solutions. In addition, it au-174 tomatically detects search stagnation and self-adapts its search strategies to escape local 175 optima (Hadka and Reed 2012; Hadka and Reed 2013). 176

In this paper, we first contribute a complete formalization of the evolutionary multiobjective direct policy search (EMODPS) approach to design closed-loop Pareto approximate operating policies for multi-purpose water reservoirs by combining DPS, nonlinear approximating networks, and the Borg MOEA. Secondly, we propose a novel EMODPS diagnostic framework to comparatively analyze the effectiveness and reliability of different policy approximation schemes (i.e., ANNs and RBFs), in order to provide practical recommendations on their use in water reservoir operating problems independently from any case-study specific calibration of the policy design process (e.g., preconditioning the decision space, tuning the optimization algorithm). Finally, we systematically review the main limitations of DP family methods in contrast to using the EMODPS approach for understanding the multi-objective tradeoffs when evaluating alternative operating policies.

The Hoa Binh water reservoir system (Vietnam) is used to demonstrate our framework. 188 The Hoa Binh is a multi-purpose reservoir that regulates the flows in the Da River, the main 189 tributary of the Red River, and is mainly operated for hydropower production and flood 190 control in Hanoi. This case study represents a relatively simple problem which, in principle, 191 should maximize the potential of SDP. As a consequence, if EMODPS met or exceeded the 192 SDP performance, we can expect that the general value of the proposed EMODPS approach 193 would increase when transitioning to more complex problems. The rest of the paper is 194 organized as follows: the next section defines the methodology, followed by the description 195 of the Hoa Binh case study. Results are then reported, while final remarks, along with issues 196 for further research, are presented in the last section. 197

#### 198 METHODS AND TOOLS

In this section, we first introduce the traditional formulation of the operating policy design problem adopted in the DP family methods and contrast it with the EMODPS formulation. The EMODPS framework has three main components: (*i*) direct policy search, (*ii*) nonlinear approximating networks, and (*iii*) multi-objective evolutionary algorithms. This section concludes with a description of the diagnostic framework used to distinguish the relative performance of ANN and RBF implementations of the proposed EMODPS approach.

#### 205 Stochastic Dynamic Programming

Water reservoir operation problems generally require sequential decisions  $\mathbf{u}_t$  (e.g., release or pumping decisions) at discrete time instants on the basis of the current system conditions

described by the state vector  $\mathbf{x}_t$  (e.g., reservoir storage). The decision vector  $\mathbf{u}_t$  is determined, 208 at each time step, by an operating policy  $\mathbf{u}_t = p(t, \mathbf{x}_t)$ . The state of the system is then altered 209 according to a transition function  $\mathbf{x}_{t+1} = f_t(\mathbf{x}_t, \mathbf{u}_t, \boldsymbol{\varepsilon}_{t+1})$ , affected by a vector of stochastic 210 external drivers  $\boldsymbol{\varepsilon}_{t+1}$  (e.g., reservoir inflows). In the adopted notation, the time subscript of a 211 variable indicates the instant when its value is deterministically known. Since SDP requires 212 that the system dynamics are known, the external drivers can only be made endogenous into 213 the SDP formulation either as state variables, described by appropriate dynamic models 214 (i.e.,  $\boldsymbol{\varepsilon}_{t+1} = f_t(\boldsymbol{\varepsilon}_t, \cdot)$ ), or as stochastic disturbances, represented by their associated pdf (i.e., 215  $\boldsymbol{\varepsilon}_{t+1} \sim \phi_t).$ 216

The combination of states and decisions over the time horizon t = 1, ..., H defines a trajectory  $\tau$ , which allows evaluating the performance of the operating policy p as follows:

$$J_p = \Psi[R(\tau)|p] \tag{1}$$

where  $R(\tau)$  defines the objective function of the problem (assumed to be a cost) and  $\Psi[\cdot]$  is a filtering criterion (e.g., the expected value) to deal with uncertainties generated by  $\varepsilon_{t+1}$ . The optimal policy  $p^*$  is hence obtained by solving the following problem:

$$p^* = \arg\min_p J_p \tag{2}$$

<sup>222</sup> subject to the dynamic constraints given by the state transition function  $\mathbf{x}_{t+1} = f_t(\mathbf{x}_t, \mathbf{u}_t, \boldsymbol{\varepsilon}_{t+1})$ . <sup>223</sup> The DP family methods solve Problem (2) by estimating the expected long-term cost of <sup>224</sup> a policy for each state  $\mathbf{x}_t$  at time t by means of the value function

$$Q_t(\mathbf{x}_t, \mathbf{u}_t) = \mathbb{E}_{\boldsymbol{\varepsilon}_{t+1}}[g_t(\mathbf{x}_t, \mathbf{u}_t, \boldsymbol{\varepsilon}_{t+1}) + \gamma \min_{\mathbf{u}_{t+1}} Q_{t+1}(\mathbf{x}_{t+1}, \mathbf{u}_{t+1})]$$
(3)

where  $Q_t(\cdot)$  is defined over a discrete grid of states and decisions,  $g_t(\cdot)$  represents the immediate (time separable) cost function associated to the transition from state  $\mathbf{x}_t$  to state  $\mathbf{x}_{t+1}$ under the decision  $\mathbf{u}_t$ , and  $\gamma \in (0, 1]$  a discount factor. With this formulation, the expected value is the statistic used to filter the uncertainty (i.e.,  $\Psi[\cdot] = E[\cdot]$ ). The optimal policy is then derived as the one minimizing the value function, namely  $p^* = \arg \min_p Q_t(\mathbf{x}_t, \mathbf{u}_t)$ .

The computation of the value function defined in eq. (3) requires the following modeling 230 assumptions (Castelletti et al. 2012): (i) the system is modeled as a discrete automaton 231 with finite domains of state, decision, and disturbance variables, with the latter described as 232 stochastic variables with an associated pdf; (ii) the objective function must be time-separable 233 along with the problem's constraints; (*iii*) the disturbance process must be uncorrelated in 234 time. Although these assumptions might appear to be restrictive, they can be applied to the 235 majority of the water resources systems by properly enlarging the state vector dimensionality 236 (Soncini-Sessa et al. 2007). For example, a duration curve can be modeled as time-separable 237 by using an auxiliary state variable accounting for the length of time. Unfortunately, the 238 resulting computation of  $Q_t(\mathbf{x}_t, \mathbf{u}_t)$  becomes very challenging in high-dimensional state and 239 decision spaces. Let  $n_x, n_u, n_{\varepsilon}$  be the number of state, decision, and disturbance variables 240 with  $N_x, N_u, N_{\varepsilon}$  the number of elements in the associated discretized domains, the compu-241 tational complexity of SDP is proportional to  $(N_x)^{n_x} \cdot (N_u)^{n_u} \cdot (N_\varepsilon)^{n_\varepsilon}$ . 242

When the problem involves multiple objectives, the single-objective optimization must 243 be repeated for every Pareto optimal point by using different scalarization values, such as 244 changing the weights used in the convex combination of the objectives (Gass and Saaty 245 1955). The overall cost of SDP to obtain an approximation of the Pareto optimal set is 246 therefore much higher, as a linear increase in the number of objectives considered yields 247 a factorial growth of the number of sub-problems to solve (i.e., a four objective problem 248 requires to solve also 4 single-objective sub-problems, 6 two-objective sub-problems, and 249 4 three-objective sub-problems (Reed and Kollat 2013; Giuliani et al. 2014)). It follows 250 that SDP cannot be applied to water systems where the number of reservoirs as well as the 251 number of objectives increases. Finally, it is worth noting that the adoption of a convex 252 combination of the objectives allows exploring only convex tradeoff curves, with gaps in 253 correspondence to concave regions. Although concave regions can be explored by adopting 254

alternative scalarization functions, such as the  $\varepsilon$ -constraint method (Haimes et al. 1971), this approach cannot be applied in the SDP framework because it violates the requirement of time-separability.

#### 258 Direct Policy Search

Direct policy search (DPS, see Sutton et al. (2000); Rosenstein and Barto (2001)) replaces the traditional SDP policy design approach, based on the computation of the value function, with a simulation-based optimization that directly operates in the policy space. DPS is based on the parameterization of the operating policy  $p_{\theta}$  and the exploration of the parameter space  $\Theta$  to find a parameterized policy that optimizes the expected long-term cost, i.e.

$$p_{\theta}^* = \arg\min_{p_{\theta}} J_{p_{\theta}} \quad \text{s.t. } \theta \in \Theta; \quad \mathbf{x}_{t+1} = f_t(\mathbf{x}_t, \mathbf{u}_t, \boldsymbol{\varepsilon}_{t+1})$$
(4)

where the objective function  $J_{p_{\theta}}$  is defined in eq. (1). Finding  $p_{\theta}^{*}$  is equivalent to find the corresponding optimal policy parameters  $\theta^{*}$ .

As reviewed by Deisenroth et al. (2011), different DPS approaches have been proposed and they differ in the methods adopted for the generation of the system trajectories  $\tau$  used in the estimation of the objective function and for the update and evaluation of the operating policies. Among them, in order to avoid the three curses of SDP and to advance the design of operating policies for multi-purpose water reservoirs, we focus on the use of an evolutionary multi-objective direct policy search (EMODPS) approach (see Figure 1) with the following features:

• Stochastic trajectory generation: the dynamic model of the system is used as simulator for sampling the trajectories  $\tau$  used for the estimation of the objective function. In principle, given the stochasticity of water systems, the model should be simulated under an infinite number of disturbance realizations, each of infinite length, in order to estimate the value of the objective function defined in eq. (1). In practice, the expected value over the probability distribution of the disturbances can be approximated with the average value over a sufficiently long time series of disturbances' realizations, either historical or synthetically generated (Pianosi et al. 2011). An alternative is represented by the analytic computation of the system trajectories (i.e., the dynamics of the state vector probability distributions with the associated decisions). However, this latter is computationally more expensive than sampling the trajectories from the system simulation, even though it can be advantageous for the subsequent policy update, as it allows the analytic computation of the gradients.

Episode-based exploration and evaluation: the quality of an operating policy  $p_{\theta}$  (and 286 of its parameter vector  $\theta$ ) is evaluated as the expected return computed on the whole 287 episode (i.e., a system simulation from t = 0, ..., H) to allow considering non-time 288 separable objectives and constraints (e.g., flow duration curves) without augmenting 289 the state vector's dimensionality. On the contrary, the step-based exploration and 290 evaluation assesses the quality of single state-decision pairs by changing the param-291 eters  $\theta$  at each time step. As in other on-line approaches, such as traditional model 292 predictive control (Bertsekas 2005), this approach requires setting a penalty function 293 on the final state (condition) of the system to account for future costs (Mayne et al. 294 2000). Yet, the definition of this penalty function requires the evaluation of the value 295 function and, hence, suffers the same limitation of DP family methods. 296

• Multi-objective: although most of DPS approaches looks at a single measure of policy performance, optimized via single-objective gradient-based optimization methods (Peters and Schaal 2008; Sehnke et al. 2010), we replace the single-objective formulation (eqs. 1-4) with a multi-objective one, where  $J_{p_{\theta}}$  and  $p_{\theta}$  represent the objective and policy vectors, respectively, that can be solved via multi-objective evolutionary algorithms (MOEAs).

The core components of the EMODPS framework have been selected to alleviate the restrictions posed by the three main curses of SDP: (*i*) EMODPS overcomes the curse of dimensionality, as it avoids the computation of the value function  $Q(\mathbf{x}_t, \mathbf{u}_t)$  (see eq. (3))

for each combination of the discretized state and decision variables, along with the biases 306 introduced by the discretization of the state, decision, and disturbance domains (Baxter et al. 307 2001). In addition, episode-based methods are not restricted to time-separable cost functions, 308 which can depend on the entire simulated trajectory  $\tau$ . (ii) EMODPS overcomes the curse of 309 modeling, as it can be combined with any simulation model as well as it can directly employ 310 exogenous information (e.g., observed or predicted inflows and precipitation) to condition 311 the decisions, without presuming either an explicit dynamic model or the estimation of any 312 pdf. (*iii*) EMODPS overcomes the curse of multiple objectives, as the combination of DPS 313 and MOEAs allows users to explore the Pareto approximate tradeoffs for up to ten objectives 314 in a single run of the algorithm (Kasprzyk et al. 2009; Reed and Kollat 2013; Giuliani et al. 315 2014). 316

Beyond these practical advantages, the general application of EMODPS does not provide 317 theoretical guarantees on the optimality of the resulting operating policies, which are strongly 318 dependent on the choice of the class of functions to which they belong and on the ability of 319 the optimization algorithm to deal with non-linear models and objectives functions, complex 320 and highly constrained decision spaces, and multiple competing objectives. Some guarantees 321 of convergence and the associated approximation bounds with respect to a known optimal 322 solution have been defined for some classes of single-objective problems, characterized by 323 time-separable and regular cost functions that can be solved with gradient-based methods 324 (Zoppoli et al. 2002; Gaggero et al. 2014). Nonetheless, EMODPS can also be employed 325 in multi-objective applications where a reference optimal solution cannot be computed due 326 to the problem's complexity, facilitating potentially good approximations of the unknown 327 optimum for a broader class of problems. 328

#### 329 Nonlinear approximating networks

The definition of a parameterized operating policy provides a mapping between the decisions  $\mathbf{u}_t$  and the policy inputs  $\mathcal{I}_t$ , namely  $\mathbf{u}_t = p_{\theta}(\mathcal{I}_t)$ . In the literature, a number of parameterizations of water reservoir operating rules have been proposed, defining the re-

lease decision as a function of the reservoir storage (Lund and Guzman 1999; Celeste and 333 Billib 2009). However, most of these rules have been derived from empirical considerations 334 and for single-objective problems, such as the design of hedging rules for flood management 335 (Tu et al. 2003) or of water supply operations (Momtahen and Dariane 2007). Indeed, if 336 prior knowledge about a (near-)optimal policy is available, an ad-hoc policy parameteriza-337 tion can be designed: parameterizations that are linear in the state variables can be used 338 when it is known that a (near-)optimal policy is a linear state feedback. However, when the 339 complexity of the system increases, more flexible structures depending on a high number of 340 parameters are required to avoid restricting the search for the optimal policy to a subspace 341 of the decision space that does not include the optimal solution. In addition, the presence 342 of multiple objectives may require to condition the decisions not only on the reservoir stor-343 age, but also on additional information (e.g., inflows, temperature, precipitation, snow water 344 equivalent (Hejazi and Cai 2009)). Two alternative approaches are available to this end: (i)345 identify a dynamic model describing each additional information and use the states of these 346 models to condition the operating policies in a DP framework (Tejada-Guibert et al. 1995; 347 Desreumaux et al. 2014); (ii) adopt approximate dynamic programming methods allowing 348 the direct, model-free use of information in conditioning the operating policies (Faber and 349 Stedinger 2001; Castelletti et al. 2010). 350

In order to ensure flexibility to the operating policy structure and to potentially condition 351 the decisions on several variables, we define the parameterized operating policy  $p_{\theta}$  by means 352 of two nonlinear approximating networks, namely Artificial Neural Networks and Gaussian 353 Radial Basis Functions. These nonlinear approximating networks have been proven to be 354 universal approximators (for a review see Tikk et al. (2003) and references therein): under 355 very mild assumptions on the activation functions used in the hidden layer, it has been shown 356 that any continuous function defined on a closed and bounded set can be approximated by 357 a three-layered ANNs (Cybenko 1989; Funahashi 1989; Hornik et al. 1989) as well as by 358 a three-layered RBFs (Park and Sandberg 1991; Mhaskar and Micchelli 1992; Chen and 359

Chen 1995). Since these features guarantee high flexibility to the shape of the parameterized function, ultimately allowing to get closer to the unknown optimum, ANNs and RBFs have become widely adopted as universal approximators in many applications (Maier and Dandy 2000; Buhmann 2003; de Rigo et al. 2005; Castelletti et al. 2007; Busoniu et al. 2011).

#### 364 Artificial Neural Networks

Using ANNs to parameterize the policy, the k-th component in the decision vector  $\mathbf{u}_t$ (with  $k = 1, ..., n_u$ ) is defined as:

$$u_t^k = a_k + \sum_{i=1}^N b_{i,k} \psi_i (\mathcal{I}_t \cdot \mathbf{c}_{i,k} + d_{i,k})$$
(5)

where N is the number of neurons with activation function  $\psi(\cdot)$  (i.e., hyperbolic tangent sigmoid function),  $\mathcal{I}_t \in \mathbb{R}^M$  the policy inputs vector, and  $a_k, b_{i,k}, d_{i,k} \in \mathbb{R}$ ,  $\mathbf{c}_{i,k} \in \mathbb{R}^M$  the ANNs parameters. To guarantee flexibility to the ANN structure, the domain of the ANN parameters is defined as  $-10,000 < a_k, b_{i,k}, \mathbf{c}_{i,k}, d_{i,k} < 10,000$  (Castelletti et al. 2013). The parameter vector  $\theta$  is therefore defined as  $\theta = [a_k, b_{i,k}, \mathbf{c}_{i,k}, d_{i,k}]$ , with  $i = 1, \ldots, N$  and  $k = 1, \ldots, n_u$ , and belongs to  $\mathbb{R}^{n_{\theta}}$ , where  $n_{\theta} = n_u(N(M+2)+1)$ .

#### 373 Radial Basis Functions

In the case of using RBFs to parameterize the policy, the k-th decision variable in the vector  $\mathbf{u}_t$  (with  $k = 1, ..., n_u$ ) is defined as:

$$u_t^k = \sum_{i=1}^N w_{i,k} \varphi_i(\mathcal{I}_t) \tag{6}$$

where N is the number of RBFs  $\varphi(\cdot)$  and  $w_{i,k}$  the weight of the *i*-th RBF. The weights are formulated such that they sum to one (i.e.,  $\sum_{i=1}^{N} w_{i,k} = 1$ ) and are non-negative (i.e.,  $w_{i,k} \ge 0 \quad \forall i, k$ ). The single RBF is defined as follows:

$$\varphi_i(\mathcal{I}_t) = \exp\left[-\sum_{j=1}^M \frac{((\mathcal{I}_t)_j - c_{j,i})^2}{b_{j,i}^2}\right]$$
(7)

where M is the number of policy inputs  $\mathcal{I}_t$  and  $\mathbf{c}_i, \mathbf{b}_i$  are the M-dimensional center and radius vectors of the *i*-th RBF, respectively. The centers of the RBF must lie within the bounded input space and the radii must strictly be positive (i.e., using normalized variables,  $\mathbf{c}_i \in [-1, 1]$  and  $\mathbf{b}_i \in (0, 1]$ , (Busoniu et al. 2011)). The parameter vector  $\theta$  is therefore defined as  $\theta = [c_{i,j}, b_{i,j}, w_{i,k}]$ , with  $i = 1, \ldots, N, j = 1, \ldots, M, k = 1, \ldots, n_u$ , and belongs to  $\mathbb{R}^{n_{\theta}}$ , where  $n_{\theta} = N(2M + n_u)$ .

#### <sup>385</sup> Multi-objective evolutionary algorithms

Multi-objective evolutionary algorithms (MOEAs) are iterative search algorithms that 386 evolve a Pareto-approximate set of solutions by mimicking the randomized mating, selec-387 tion, and mutation operations that occur in nature (Deb 2001; Coello Coello et al. 2007). 388 These mechanisms allow MOEAs to deal with challenging multi-objective problems char-389 acterized by multi-modality, nonlinearity, stochasticity and discreteness, thus representing 390 a promising alternative to gradient-based optimization methods in solving multi-objective 391 water reservoirs problems (see Nicklow et al. (2010) and Maier et al. (2014) and references 392 therein). 393

In this paper, we use the self-adaptive Borg MOEA (Hadka and Reed 2013), which 394 employs multiple search operators that are adaptively selected during the optimization, based 395 on their demonstrated probability of generating quality solutions. The Borg MOEA has been 396 shown to be highly robust across a diverse suite of challenging multi-objective problems, 397 where it met or exceeded the performance of other state-of-the-art MOEAs (Hadka and 398 Reed 2012; Reed et al. 2013). In addition to adaptive operator selection, the Borg MOEA 399 assimilates several other recent advances in the field of MOEAs, including an  $\varepsilon$ -dominance 400 archiving with internal algorithmic operators to detect search stagnation, and randomized 401 restarts to escape local optima. The flexibility of the Borg MOEA to adapt to challenging, 402 diverse problems makes it particularly useful for addressing EMODPS problems, where the 403 shape of the operating rule and its parameter values are problem-specific and completely 404 unknown a priori. 405

#### 406 Diagnostic framework

In this work, we apply a diagnostic framework developed from the one in Hadka and 407 Reed (2012) to comparatively analyze the potential of the ANN and RBF policy parameter-408 izations in solving EMODPS problems with no specific tuning of the policy design process. 400 Since the presence of multiple objectives does not yield a unique optimal solution, but a 410 set of Pareto optimal solutions, assessing the effectiveness of the policy design results (i.e., 411 how close the solutions found are to the optimal ones) requires to evaluate multiple metrics, 412 such as the distance of the final solutions from the Pareto optimal front or its best known 413 approximation (i.e., reference set), the coverage of the non-dominated space, and the extent 414 of the non-dominated front (Maier et al. 2014). In this work, we adopt three formal met-415 rics, namely generational distance, additive  $\varepsilon$ -indicator, and hypervolume indicator, which 416 respectively account for convergence, consistency, and diversity (Knowles and Corne 2002; 417 Zitzler et al. 2003). In addition, due to the stochastic nature of the evolutionary algorithms 418 (which can be affected by random effects in initial populations and runtime search oper-419 ators), each optimization was run for multiple random generator seeds. The reliability of 420 the ANN and RBF policy search is evaluated as the probability of finding a solution that 421 is better or equal to a certain performance threshold in a single run, which measures the 422 variability in the solutions' effectiveness for repeated optimization trials. 423

424

The generational distance  $I_{GD}$  measures the average Euclidean distance between the points in an approximation set S and the nearest corresponding points in the reference set  $\bar{S}$ , and it is defined as

$$I_{GD}(S,\bar{S}) = \frac{\sqrt{\sum_{\mathbf{s}\in S} d_{\mathbf{s}}^2}}{n_S}$$
(8a)

428 with

$$d_{\mathbf{s}} = \min_{\bar{\mathbf{s}}\in\bar{S}} \sqrt{\sum_{i=1}^{k} [J^i(\mathbf{s}) - J^i(\bar{\mathbf{s}})]^2}$$
(8b)

where  $n_S$  is the number of points in S, and  $d_s$  the minimum Euclidean distance between each point in S and  $\overline{S}$ .  $I_{GD}$  is a pure measure of convergence and the easiest to satisfy, requiring only a single solution close to the reference set to attain ideal performance.

<sup>432</sup> The additive  $\varepsilon$ -indicator  $I_{\varepsilon}$  measures the worst case distance required to translate an <sup>433</sup> approximation set solution to dominate its nearest neighbour in the reference set, defined as

$$I_{\varepsilon}(S,\bar{S}) = \max_{\bar{\mathbf{s}}\in\bar{S}} \min_{\mathbf{s}\in S} \max_{1\leq i\leq k} (J^{i}(\mathbf{s}) - J^{i}(\bar{\mathbf{s}}))$$
(9)

This metric is very sensitive to gaps in tradeoff and can be viewed as a measure of an approximation set's consistency with the reference set, meaning that all portions of the tradeoff are present (Hadka and Reed 2012). Additionally, it captures diversity because of its focus on the worst case distance. If a Pareto approximate set S has gaps, then solutions from other regions must be translated much further distances to dominate its nearest neighbour in the reference set  $\bar{S}$ , dramatically increasing the  $I_{\varepsilon}$  value.

Finally, the hypervolume measures the volume of objective space dominated by an approximation set, capturing both convergence and diversity. It is the most challenging of the three metrics to satisfy. The hypervolume indicator  $I_H$  is calculated as the difference in hypervolume between the reference set  $\bar{S}$ , and an approximation set S, defined as

$$I_H(S,\bar{S}) = \frac{\int \alpha_S(\mathbf{s})ds}{\int \alpha_{\bar{S}}(\bar{\mathbf{s}})d\bar{s}}$$
(10a)

444 with

$$\alpha(\mathbf{s}) = \begin{cases} 1 & \text{if } \exists \mathbf{s}' \in S \text{ such that } \mathbf{s}' \preceq \mathbf{s} \\ 0 & \text{otherwise} \end{cases}$$
(10b)

Overall, a good set of Pareto approximate policies is characterized by low values of the first two metrics and a high value of the third one.

#### 447 CASE STUDY DESCRIPTION

The Hoa Binh is a multi-purpose regulated reservoir in the Red River basin, Vietnam

(Figure 2). The Red River drains a catchment of  $169,000 \text{ km}^2$  shared by China (48%), 449 Vietnam (51%), and Laos (1%). Among the three main tributaries (i.e., Da, Thao, and Lo 450 rivers), the Da River is the most important water source, contributing for 42% of the total 451 discharge at Hanoi. Since 1989, the discharge from the Da River has been regulated by the 452 operation of the Hoa Binh reservoir, which is one of the largest water reservoirs in Vietnam, 453 characterized by a surface area of about  $198 \text{ km}^2$  and an active storage capacity of about 454 6 billion  $m^3$ . The dam is connected to a power plant equipped with eight turbines, for a 455 total design capacity of 1920 MW, which guarantees a large share of the national electricity 456 production. Given the large storage capacity, the operation of Hoa Binh has also a key role 457 for flood mitigation in Hanoi in the downstream part of the Red River catchment (Castelletti 458 et al. 2012). In recent years, other reservoirs have been constructed on both the Da and Lo 459 rivers (see the yellow triangles in Figure 2). However, given the limited data available since 460 these reservoirs have started operating, they are not considered in this work. 461

#### <sup>462</sup> Model and objectives formulation

The system is modeled by a combination of conceptual and data-driven models assuming and a modeling and decision-making time-step of 24 hours. The Hoa Binh dynamics is described by the mass balance equation of the water volume  $s_t^{HB}$  stored in the reservoir, i.e.

$$s_{t+1}^{HB} = s_t^{HB} + q_{t+1}^D - r_{t+1}$$
(11)

where  $q_{t+1}^D$  is the net inflow to the reservoir in the interval [t, t + 1) (i.e., inflow minus evaporation losses) and  $r_{t+1}$  is the volume released in the same interval. The release is defined as  $r_{t+1} = f(s_t^{HB}, u_t, q_{t+1}^D)$ , where  $f(\cdot)$  describes the nonlinear, stochastic relation between the decision  $u_t$ , and the actual release  $r_{t+1}$  (Piccardi and Soncini-Sessa 1991). The flow routing from the reservoir to the city of Hanoi is instead described by a data-driven feedforward neural network, providing the level in Hanoi given the Hoa Binh release  $(r_{t+1})$  and the Thao  $(q_{t+1}^T)$  and Lo  $(q_{t+1}^L)$  discharges. The description of the Hoa Binh net inflows  $(q_{t+1}^D)$  and the

flows in the Thao  $(q_{t+1}^T)$  and Lo  $(q_{t+1}^L)$  rivers depends on the approach adopted: with SDP, 473 they are modeled as stochastic disturbances; with EMODPS, they are not explicitly modeled 474 as this approach allows to directly embed exogenous information into the operating policies. 475 Further details about the model of the Hoa Binh system can be found in Castelletti et al. 476 (2012) and Castelletti et al. (2013). 477

The two conflicting interests affected by the Hoa Binh operation are modeled using the 478 following objective formulations, evaluated over the simulation horizon H: 479

- Hydropower production  $(J^{hyd})$ : the daily average hydropower production (kWh/day) 480 at the Hoa Binh hydropower plant, to be maximized, defined as 481

$$J^{hyd} = \frac{1}{H} \sum_{t=0}^{H-1} H P_{t+1}$$
(12)  
with  $HP_{t+1} = \left(\eta g \gamma_w \bar{h}_t q_{t+1}^{Turb}\right) \cdot 10^{-6}$ 

where  $\eta$  is the turbine efficiency,  $g = 9.81 \ (m/s^2)$  the gravitational acceleration, 482  $\gamma_w = 1000 \ (\text{kg/m}^3)$  the water density,  $\bar{h}_t$  (m) the net hydraulic head (i.e., reservoir 483 level minus tailwater level),  $q_{t+1}^{Turb}$  (m<sup>3</sup>/s) the turbined flow; 484

Flooding  $(J^{flo})$ : the daily average excess level  $h_{t+1}^{Hanoi}(\text{cm}^2/\text{day})$  in Hanoi with respect 485 to the flooding threshold  $\bar{h} = 950$  cm, to be minimized, defined as 486

$$J^{flo} = \frac{1}{H} \sum_{t=0}^{H-1} max (h_{t+1}^{Hanoi} - \bar{h}, 0)^2$$
(13)

487

where  $h_{t+1}^{Hanoi}$  is the level in Hanoi estimated by the flow routing model, which depends on the Hoa Binh release  $(r_{t+1})$  along with the Thao  $(q_{t+1}^T)$  and Lo  $(q_{t+1}^L)$  discharges. 488

It is worth noting that the proposed model and objective formulations are defined as 489 Markov Decision Processes (Soncini-Sessa et al. 2007) to allow comparing the results of 490 EMODPS with traditional DP-based solutions. A more realistic representation would re-491

quire the development of hydrological models describing the rivers catchments and the use 492 of a flooding objective function that is not time-separable (e.g., the duration of the flood 493 event, which may induce dykes breaks when exceeding critical thresholds). Yet, these al-494 ternatives would enlarge the state vector dimensionality beyond the SDP limits. Moreover, 495 the curse of multiple objectives narrows the number of water-related interests that can be 496 considered, preventing a better understanding of the full set of tradeoffs (e.g., flood peaks 497 vs flood duration) and ignoring less critical sectors (e.g., irrigation and environment). The 498 adopted formulations therefore represent a relatively simplified system representation which, 499 in principle, should maximize the potential of SDP. Given the heuristic nature of EMODPS, 500 which has no guarantee of optimality, we use SDP as a benchmark to evaluate the quality 501 of the approximation attained by the EMODPS operating policies. If EMODPS met or ex-502 ceeded the SDP performance, the general value of the proposed EMODPS approach would 503 increase by including additional model/objective complexities. 504

#### 505 Computational Experiment

The Hoa Binh operating policies are parameterized by means of three-layered nonlinear 506 approximating networks, where different numbers of neurons and basis functions are tested. 507 According to Bertsekas (1976), the minimum set of policy inputs required to produce the 508 best possible performance is the state of the system  $\mathbf{x}_t$ , possibly coupled with a time index 509 (e.g., the day of the year) to take into account the time-dependency and cyclostationarity of 510 the system and, consequently, of the operating policy. However, according to previous works 511 (Pianosi et al. 2011; Giuliani et al. 2014), the operating policy of the Hoa Binh reservoir 512 benefits from the consideration of additional variables, which cannot be employed in DP 513 methods without enlarging the state-vector dimensionality. In particular, the best operation 514 of the Hoa Binh reservoir is obtained by conditioning the operating policies upon the following 515 input variables  $\mathcal{I}_t = [sin(2\pi t)/365, cos(2\pi t)/365, s_t^{HB}, q_t^D, q_t^{lat}]$ , where  $q_t^{lat} = q_t^T + q_t^L$  is the 516 lateral inflow accounting for the Thao and Lo discharges. The role of the previous day inflow 517 observations  $q_t^D$  and  $q_t^{lat}$  is key in enlarging the information on the current system condition, 518

particularly with respect to the flooding objective in Hanoi, which depends on both the Hoa
Binh releases as well as on the lateral flows of Thao and Lo rivers.

The EMODPS optimization of the parameterized operating policies employs the Borg 521 MOEA. Since it has been demonstrated to be relatively insensitive to the choice of param-522 eters, we use the default algorithm parameterization suggested by Hadka and Reed (2013). 523 Epsilon-dominance values equal to 5000 for  $J^{hyd}$  and 5 for  $J^{flo}$  are used to set the resolution 524 of the two operating objectives. Each optimization was run for 500,000 function evalua-525 tions. To improve solution diversity and avoid dependence on randomness, the solution set 526 from each formulation is the result of 20 random optimization trials. In the analysis of the 527 runtime search dynamics, the number of function evaluations (NFE) was extended to 2 mil-528 lions. Each optimization was run over the horizon 1962-1969, which has been selected as it 529 comprises normal, wet, and dry years. The final set of Pareto approximate policies for each 530 policy structure is defined as the set of non-dominated solutions from the results of all the 20 531 optimization trials. The three metrics (i.e., generational distance, additive  $\varepsilon$ -indicator, and 532 hypervolume indicator) are computed with respect to the overall best known approximation 533 of the Pareto front, obtained as the set of non-dominated solutions from the results of all the 534 280 optimization runs (i.e., 2 approximating networks times 7 structures times 20 seeds). In 535 total, the comparative analysis comprises 220 million simulations and requires approximately 536 1,220 computing hours on an Intel Xeon E5-2660 2.20 GHz with 32 processing cores and 96 537 GB Ram. However, it should be noted that our computational experiment is more rigorous 538 than would be necessary in practice and it was performed to support a rigorous diagnostic 539 assessment of the ANN and RBF policy parameterizations. The EMODPS policy design 540 reliably attained very high fidelity approximations of the Pareto front in each optimization 541 run with approximately 150,000 NFE, corresponding to only 50 computing minutes. 542

The SDP solutions were designed by computing the value function (eq. 3) over the 2dimensional state vector  $\mathbf{x}_t = [t, s_t]$  and the Hoa Binh release decision  $u_t$ . The two objectives are aggregated trough a convex combination as the  $\varepsilon$ -constraint method would violate the SDP requirement of time-separability. The policy performance are then evaluated via simulation of the same model used in the EMODPS experiments. The stochastic external drivers
are represented as follows:

$$q_{t+1}^D \sim \mathcal{L}_t$$

$$q_{t+1}^T = \alpha^T q_{t+1}^D + \varepsilon_{t+1}^T$$

$$q_{t+1}^L = \alpha^L q_{t+1}^D + \varepsilon_{t+1}^L$$
(14)

where  $\mathcal{L}_t$  is a log-normal probability distribution and the coefficients  $(\alpha^T, \alpha^L)$  describe the spatial correlation of the inflow processes, with normally distributed residuals  $\varepsilon_{t+1}^T \sim \mathcal{N}^T$  and  $\varepsilon_{t+1}^L \sim \mathcal{N}^L$ . The models of the inflows, namely the three probability distributions  $\mathcal{L}_t, \mathcal{N}^T, \mathcal{N}^L$ as well as the coefficients  $(\alpha^T, \alpha^L)$ , were calibrated over the horizon 1962-1969 to provide the same information employed in the EMODPS approach.

The SDP problem formulation hence comprises two state variables, one decision variable, 554 and three stochastic disturbances. Preliminary experiments allow calibrating the discretiza-555 tion of state, decision, and disturbance vectors as well as the number of weights combinations 556 for aggregating the two competing objectives to identify a compromise between modeling 557 accuracy and computational requirements. Each solution designed via SDP required around 558 45 computing minutes. In order to obtain an equivalent exploration of the Pareto front as 559 in the EMODPS approach, in principle the SDP should be run for 40 different combinations 560 of the objectives, corresponding to 30 computing hours. Yet, the non-linear relationships 561 between the values of the weights and the corresponding objectives value does not guaran-562 tee to obtain 40 different solutions as most of them are likely to be equivalent or Pareto 563 dominated. Despite a very accurate tuning of the objectives' weights, we obtained only four 564 Pareto approximate solutions. Finally, the cost of developing the inflows models should be 565 also considered in the estimation of the overall effort required by the SDP, whereas in the 566 EMODPS case such cost is null given the possibility of directly employing the exogenous 567 information. 568

#### 569 **RESULTS**

In this section, we first use our EMODPS diagnostic framework to identify the most effective and reliable policy approximation scheme for the Hoa Binh water reservoir problem. Secondly, we validate the EMODPS Pareto approximate policies by contrasting them with SDP-based solutions. Finally, we analyze one potentially interesting compromise solution to provide effective recommendation supporting the operation of the Hoa Binh reservoir.

#### <sup>575</sup> Identification of the operating policy structure

The first step of the EMODPS diagnostic framework aims to identify the best parameter-576 ized operating policy's structure in terms of number of neurons (for ANN policies) or basis 577 functions (for RBF policies), for a given number M = 5 of policy input variables. Figure 3 578 shows the results for seven different policy structures with the number of neurons and basis 579 functions increasing from n = 4 to n = 16. The performance of the resulting Pareto approx-580 imate operating policies, computed over the optimization horizon 1962-1969, are illustrated 581 in Figure 3a, with the arrows identifying the direction of preference for each objective. The 582 ideal solution would be a point in the top-left corner of the figure. The figure shows the 583 reference set identified for each policy structure, obtained as the set of non-dominated so-584 lutions across the 20 optimization trials performed. The overall reference set, obtained as 585 the set of non-dominated solutions from the results of all the 280 optimization runs (i.e., 2) 586 approximating networks times 7 structures times 20 seeds), is represented by a black dotted 587 line. Comparison of the best Pareto approximate sets attained across all random seed trials 588 changing the structures of both ANNs and RBFs, namely the Pareto approximate solutions 589 represented by different shapes, does not show a clear trend of policy performance improve-590 ment with increasing numbers of neurons or basis functions. The results in Figure 3a attest a 591 general superiority of the RBF policies over the ANN ones, particularly in the exploration of 592 the tradeoff region with the maximum curvature of the Pareto front (i.e., for  $J^{flo}$  values be-593 tween 100 and 200, RBFs allows attaining higher hydropower production). The ANN policies 594 do outperform the RBF ones in terms of maximum hydropower production, although this 595

small difference is concentrated in a restricted range of  $J^{flo}$ , and, likely, not decision-relevant.

597

In order to better analyze the effectiveness and the reliability in attaining good approx-598 imations of the Pareto optimal set using different ANN/RBF structures, we computed the 599 three metrics of our diagnostic framework on the solutions obtained in each optimization 600 run. The metrics are evaluated with respect to the best known approximation of the Pareto 601 front, namely the overall reference set (i.e., the black dotted line in Figures 3a). Figures 602 3b-d report the best (solid bars) and average (transparent bars) performance in terms of 603 generational distance  $I_{GD}$ , additive  $\varepsilon$ -indicator  $I_{\varepsilon}$ , and hypervolume indicator  $I_H$ , respec-604 tively. Effective policy parameterizations are characterized by low values of  $I_{GD}$  and  $I_{\varepsilon}$ , and 605 high values of  $I_{H}$ . The deviations between the best and the average metric values reflect 606 the reliability of the policy design, with large deviations identifying low reliable structures. 607 In contrast with the results in Figure 3a, the values of the metrics show substantial dif-608 ferences between ANNs and RBFs as well as their dependency on the number of neurons 609 and basis functions. The average metrics of RBF policies are consistently better than the 610 ones of ANN policies. Moreover, the average performance of ANN policies degrade when 611 the number of neurons increases (except for n = 4, where the number of ANN inputs is 612 larger than the number of neurons) probably because ANNs are overfitting the data, while 613 the RBF policies seem to be less sensitive to the number of basis. It is worth noting that the 614 gap between RBFs and ANNs decreases when looking at the best optimization run. This 615 result suggests that the ANN policy parameterization is very sensitive to the initialization 616 and the sequence of random operators employed during the Borg MOEA search, probably 617 due to the larger domain of the ANN parameters with respect to the RBF ones. In the case 618 of RBFs, indeed, the parameter space is the Cartesian product of the subsets [-1,1] for 619 each center  $c_{j,i}$  and (0,1] for each radius  $b_{j,i}$  and weight  $w_{i,k}$ . In the case of ANNs, instead, 620 parameters have no direct relationship with the policy inputs. In this work, the domain 621  $-10,000 < a_k, b_{i,k}, \mathbf{c}_{i,k}, d_{i,k} < 10,000$  is used as in Castelletti et al. (2013) to guarantee flex-622

<sup>623</sup> ibility to the ANN structure and prevents that any Pareto approximate solution is excluded<sup>624</sup> a priori.

625

To further compare the performance of RBFs and ANNs, in the second step of the 626 analysis we perform a more detailed assessment of the reliability of attaining high quality 627 Pareto approximations for alternative operating policy structures. To this purpose, we define 628 the reliability of the ANN and RBF policy search as the probability of finding a solution 629 that is better or equal to a certain performance threshold (i.e., 75% or 95%) in a single 630 optimization run, which measures the variability in the solutions' effectiveness for repeated 631 optimization trials. Figure 4 illustrates the probability of attainment with a 75% (panel a) 632 and 95% (panel b) threshold, along with a representative example of these thresholds in the 633 objective space (panel c). Figure 4a shows that the ANN policies are not able to consistently 634 meet the 75% threshold, even in terms of  $I_{GD}$  which is generally considered the easiest metric 635 to meet requiring only a single solution close to the reference set. As shown in Figure 4c, not 636 attaining 75% in  $I_{GD}$  means to have a very poor understanding of the 2-objective tradeoff, 637 with almost no information on the left half of the Pareto front. The thresholds on  $I_{\varepsilon}$  are 638 instead fairly strict, as this metric strongly penalizes the distance from the knee region of 639 the reference set. The results in Figure 4a demonstrates the superiority of the RBF policy 640 parameterizations, which attain 75% of the best metric value with a reliability of 100%641 independently from the number of basis functions. Assuming that the 75% approximation 642 can be an acceptable approximation level of the Pareto optimal set, these results imply that 643 the Hoa Binh policy design problem can likely be solved by a single optimization run with 644 an RBF policy. However, Figure 4b shows that if the 95% level was required, it would be 645 necessary to run multiple random seeds and to accumulate the best solutions across them. 646

The results in Figure 4 also allow the identification of the most reliable structure of the operating policies in terms of number of neurons and basis functions. Results in Figure 4a show that the most reliable ANN policy relies on 6 neurons, which attains the highest

reliability in  $I_{\varepsilon}$  and  $I_{H}$ , while all the RBF policies are equally reliable. By considering a 650 stricter threshold (i.e., 95%), results in Figure 4b show that the most reliable RBF policy, 651 particularly in terms of convergence and diversity (i.e., hypervolume indicator), requires 6 652 or 8 basis functions. Note that attaining 95% in terms of  $I_{\varepsilon}$  resulted to be particularly chal-653 lenging (i.e., probabilities around 10-15%) and, as illustrated in Figure 4c this threshold is 654 almost equivalent to require the identification of the best known approximation of the Pareto 655 front in a single run. In the following, we select the 6-basis structure because it depends on 656 a lower number of parameters and allows a better comparison with the 6 neurons ANNs. 657

658

The last step of the analysis looks at the runtime evolution of the Borg MOEA search 659 to ensure that the algorithm's search is at convergence. To this purpose, we run a longer 660 optimization with 2 millions function evaluations for a 6 neurons ANN policy and a 6 basis 661 RBF policy, with 20 optimization trials for each approximating network. In each run, we 662 track the search progress by computing the values of  $I_{GD}$ ,  $I_{\varepsilon}$ , and  $I_{H}$  every 1,000 function 663 evaluations until the first 50,000 evaluations and, then, every 50,000 until 2 millions. The 664 runtime search performance are reported in Figure 5 as a function of the number of function 665 evaluations used. The values of  $I_{GD}$  in Figure 5a show that few function evaluations (i.e., 666 around 250,000) allows the identification of solutions close to the reference set identified from 667 the results obtained at the end of the optimization (i.e., after 2 million function evaluations) 668 across the 20 random optimization trials performed for each approximating network (i.e., 6 669 neurons ANN and 6 basis RBF). The performance in terms of  $I_{GD}$  of both ANN and RBF 670 policies are then almost equivalent from 250,000 to 2 millions function evaluations. 671

<sup>672</sup> A higher number of function evaluations is instead necessary to reach full convergence in <sup>673</sup> the other two metrics, namely  $I_{\varepsilon}$  and  $I_H$  illustrated in Figures 5b-c, respectively. In general, <sup>674</sup> the runtime analysis of these two metrics further confirm the superiority of the RBF operating <sup>675</sup> policies over the ANN ones, both in terms of consistency (i.e.,  $I_{\varepsilon}$ ) as well as convergence <sup>676</sup> and diversity (i.e.,  $I_H$ ). Such a superiority of RBFs is evident from the beginning of the

search, when it is probably due the larger dimensionality of the ANN parameters' domain, 677 which increases the probability of having a poor performing initial population. However, the 678 Borg MOEA successfully identifies improved solutions for both ANN and RBF policies in few 679 runs, with diminishing returns between 100,000 and 200,000 function evaluations. The search 680 progress stops around 400,000 function evaluations, with the RBF policies that consistently 681 outperform the ANN ones. The limited improvements in the performance of each solution 682 from 400,000 to 2 millions demonstrate the convergence of the Borg MOEA search for both 683 ANNs and RBFs, guaranteeing the robustness of the results previously discussed, which were 684 obtained with 500,000 functions evaluations. 685

#### <sup>686</sup> Validation of EMODPS policy performance

The performance of the operating policies discussed in the previous section is computed 687 over the optimization horizon 1962-1969. To validate this performance, the designed oper-688 ating policies are re-evaluated via simulation over a different horizon, namely 1995-2004, to 689 estimate their effectiveness under different hydroclimatic conditions. We focus the analysis 690 on the most reliable policy structures resulting from the previous section, using a 6 neurons 691 ANN and a 6 basis RBF parameterization. The comparison between the performance over 692 the optimization and the validation horizons is illustrated in Figure 6a, which reports the 693 reference set obtained in the two cases across the 20 optimization trials. It is not surprising 694 that the performance attained over the optimization horizon (transparent solutions) degrade 695 when evaluated over the validation horizon (opaque solutions) since the two sets are indepen-696 dently used in the analysis. Although both ANNs and RBFs successfully explore different 697 tradeoffs between  $J^{hyd}$  and  $J^{flo}$  over the optimization horizon, the difference in performance 698 between optimization and validation clearly demonstrate that RBF operating policies out-690 perform the ANN ones. This can be explained as a consequence of the ANNs over-fitting 700 during the optimization. Indeed, although a subset of ANN policies is Pareto dominating 701 some RBF solutions over the optimization horizon (i.e., for  $J^{flo}$  values between 220 and 702 300), the ANN Pareto approximate front is completely dominated in validation by the RBF 703

solutions. The designed ANN policies seem to be over fit on the hydroclimatic conditions on 704 which they were trained and suffering from too much parametric complexity. Consequently, 705 the ANN policies fail to manage unforeseen situations. Conversely RBFs maintains good 706 performance over the validation horizon, with the corresponding Pareto front that presents 707 less gaps and with a more consistent exploration of the tradeoff between the two objectives. 708 Figure 6b contrasts the performance of the RBF policies with solutions designed via 709 Stochastic Dynamic Programming (represented by black circles) over the validation horizon 710 1995-2004. To provide a fair comparison, we illustrate both the RBF solutions conditioned 711 upon  $\mathcal{I}_t = [sin(2\pi t)/365, cos(2\pi t)/365, s_t^{HB}, q_t^D, q_t^{lat}]$  (red crosses) and, those obtained by 712 conditioning the decisions on the same variables employed by SDP, namely the day of the 713 year t and the Hoa Binh storage  $s_t^{HB}$  (magenta crosses). Results demonstrate that, de-714 spite the theoretical guarantee of optimality, SDP solutions produce a significantly lower 715 performance than EMODPS even with basic information. The two main reasons for this 716 are that SPD uses a simplified representation of the spatial and temporal correlation of 717 the inflows and a discretization of state, decision, and disturbance domains. Optimization 718 experiments with SDP using finer discretization grids (not shown for brevity) demonstrate 719 that improvements enabled by finer resolution would be marginal. In contrast, we expect 720 that SDP performance would likely increase by improving the model of the inflows, either 721 by using an autoregressive model to characterize their autocorrelation in time or by extend-722 ing the time-series to better estimate their pdf and their spatial correlation. However, this 723 refinement would further increase the computational requirements of SDP. In addition, the 724 difficulty of balancing the two objectives when aggregated through a convex combination 725 produces multiple Pareto dominated or overlapping solutions, ultimately limiting the explo-726 ration of the tradeoff between  $J^{hyd}$  and  $J^{flo}$ . Moreover, this objectives' aggregation provides 727 a convex approximation of the Pareto front and prevents the design of solutions in concave 728 regions, resulting in large gaps among the SDP solutions. This limitation does not affect the 729 EMODPS approach, which indeed identifies Pareto approximate sets with concave region in 730

<sup>731</sup> correspondence to the gaps in the SDP solutions. Finally, the possibility of directly employ<sup>732</sup> ing exogenous information in conditioning the decisions successfully enhances the resulting
<sup>733</sup> policy performance, with the red solutions that completely dominate the magenta and black
<sup>734</sup> Ones.

#### <sup>735</sup> Analysis of the EMODPS operating policy

In order to provide effective recommendation supporting the operation of the Hoa Binh 736 reservoir, we select a potential compromise solution (see Figure 6b) and we analyze the 737 corresponding operating policy. Figure 7a provides a multivariate representation of the 738 multi-input single-output RBF policy, approximated with an ensemble of 5,000 elements 739 obtained via Latin Hypercube Sampling of the policy inputs domains. The parallel-axes plot 740 represents each release decision  $u_t$  (reported on the first axis and highlighted by the green 741 color ramp) as a line crossing the other axes at the values of the corresponding policy inputs 742 (i.e., the day of the year t, the Hoa Binh storage  $s_t$ , and the previous day flow observations 743 of the Da River  $q_t^D$  and of the lateral contribution of Thao and Lo Rivers  $q_t^{lat}$ , respectively). 744 The figure shows that the highest release decisions (dark green lines) are concentrated at the 745 beginning of the monsoon season (i.e., May and June), when it is necessary to drawdown 746 the reservoir storage to make space for the flood peak, while are less dependent on the Hoa 747 Binh storage or the flow in the Da river. As expected, since the policy under consideration 748 is a compromise between the two objectives, it ensures flood protection by suggesting high 749 releases when the flows in the Thao and Lo rivers are small. Focusing on the second axis, 750 representing the day of the year, it is possible to appreciate the cyclostationary behavior of 751 the operating policy, which provides similar release decisions (i.e., mid-tone green lines) at 752 the beginning (bottom) and at the end (top) of the year. 753

Further details are provided by Figure 7b-d, which represents the release decision projected as a function of the reservoir storage, with the colors illustrating how the release decision changes depending on the day of the year (panel b), the flow in the Da River (panels c-d), and the lateral flow in Thao and Lo Rivers (panel e). Figure 7b confirms the

cyclostationary behavior of the operating policy throughout the year (for fixed, intermediate 758 values of flow in the Da River as well as in the Thao and Lo Rivers). The release decision 759 is indeed increasing to make room for the incoming flood before and during the monsoon 760 season, from May (green lines) to August (blue lines). Then, after the monsoon, it decreases 761 and the operation at the end of the year is equivalent to the one at the beginning of the year 762 (red lines). Figure 7c shows the release decision as a function of the Hoa Binh storage on 763 January the  $1^{st}$  for different values of flow in the Da River (and a fixed intermediate value 764 of flow in the Thao and Lo Rivers). In this case, according to the value of the inflow (i.e., 765 moving from light to dark green) the release decision increases to maximize the hydropower 766 production, while maintaining a high and constant water level in the Hoa Binh reservoir. 767 Although we are considering a compromise policy, such increasing releases are acceptable 768 also in terms of flood protection because the monsoon season is far in the future. The mod-769 ification of the policy during the monsoon season is evident in Figure 7d, which shows again 770 the release decision as a function of the Hoa Binh storage for different values of flow in the 771 Da River (and a fixed intermediate value of flow in the Thao and Lo Rivers) but on May 772 the  $1^{st}$ . In this case the release decision is first increasing with the inflow but, when this 773 latter exceeds 9,000 m<sup>3</sup>/s, it starts decreasing to reduce the flood costs in Hanoi. Finally, 774 Figure 7e represents the dual situation, namely the release decision as a function of the Hoa 775 Binh storage on May the  $1^{st}$  for different values of flow in the Thao and Lo Rivers (and a 776 fixed intermediate value of flow in the Da River). In this case, effective flood protection is 777 obtained by decreasing the release decision when the lateral flow increases (i.e., moving from 778 light to dark green lines). 779

#### 780 CONCLUSIONS

The paper formalizes and demonstrates the potential of the evolutionary multi-objective direct policy search approach in advancing water reservoirs operations. The method combines direct policy search method, nonlinear approximating networks, and multi-objective evolutionary algorithms to design Pareto approximate operating policies for multi-purpose water reservoirs. The regulation of the Hoa Binh water reservoir in Vietnam is used as acase study.

The comparative analysis of two widely used nonlinear approximating networks (i.e., 787 Artificial Neural Networks and Gaussian Radial Basis Functions) for the parameterization 788 of the operating policy suggests the general superiority of RBFs over ANNs. Results show 789 that RBF solutions are more effective that ANN ones in designing Pareto approximate 790 policies for the Hoa Binh reservoir, with better performance attained by the associated 791 Pareto fronts in terms of convergence, consistency, and diversity. Moreover, the adopted 792 EMODPS diagnostic framework demonstrates that the search of RBF policies is more reliable 793 than using ANNs, thus guaranteeing an high probability of designing high quality solutions. 794 Finally, the performance of RBF policies consistently outperform the ANN ones also when 795 simulated on a different horizon with respect to the one used for the optimization. Although 796 accurate calibration and preconditioning of ANN policies have been shown to improve their 797 performance (Castelletti et al. 2013), they require a priori information about the shape of 798 the optimal policy. On the contrary, RBF operating policies successfully attain high quality 799 results without any tuning or preconditioning of the policy design process, thus representing 800 a potentially effective, case study-independent option for solving EMODPS problems. In 801 addition, although the Hoa Binh policy design problem formulation as a 2-objective Markov 802 Decision Process should maximize the potential of Stochastic Dynamic Programming, our 803 results demonstrate that EMODPS successfully improves the SDP solutions, showing the 804 potential to overcome most of the limitations of DP family methods. The general value of 805 the proposed EMODPS approach would further increase when transitioning to more complex 806 problems. Finally, the analysis of the RBF policy shows physically sound interpretations, 807 favoring its acceptability for the reservoir operators and contributing quantitative practical 808 recommendation to improve the Hoa Binh regulation. 809

Future research efforts will focus on testing the scalability of EMODPS with respect to the dimensionality of the state and decision vectors as well as to the number of objectives, particularly to support the use of EMODPS in multireservoir systems (Biglarbeigi et al. 2014), possibly including robustness criteria to face global change (Herman et al. 2015). Moreover, the scope of the comparative analysis might be enlarged by including other approximators, such as fuzzy systems or support vector machine. Finally, a diagnostic assessment on different state-of-the-art MOEAs in EMODPS problems will be developed.

#### 817 ACKNOWLEDGEMENT

This work was partially supported by the *IMRR - Integrated and sustainable water Management of the Red-Thai Binh Rivers System in changing climate* research project funded by the Italian Ministry of Foreign Affair as part of its development cooperation program. Francesca Pianosi was supported by the Natural Environment Research Council (Consortium on Risk in the Environment: Diagnostics, Integration, Benchmarking, Learning and Elicitation (CREDIBLE); grant number NE/J017450/1).

#### 824 **REFERENCES**

- Ansar, A., Flyvbjerg, B., Budzier, A., and Lunn, D. (2014). "Should we build more large dams? The actual costs of hydropower megaproject development." *Energy Policy*, 69, 43–56.
- Baxter, J., Bartlett, P., and Weaver, L. (2001). "Experiments with infinite-horizon, policygradient estimation." J. Artif. Intell. Res. (JAIR), 15, 351–381.
- Bellman, R. (1957). Dynamic programming. Princeton University Press, Princeton.
- Bertsekas, D. (1976). Dynamic programming and stochastic control. Academic Press, New
  York.
- Bertsekas, D. (2005). "Dynamic programming and suboptimal control: a survey from ADP
  to MPC." *European Journal of Control*, 11(4-5).
- Bertsekas, D. and Tsitsiklis, J. (1996). Neuro-dynamic programming. Athena Scientific, Belmont, MA.
- <sup>837</sup> Biglarbeigi, P., Giuliani, M., and Castelletti, A. (2014). "Many-objective direct policy search

- in the Dez and Karoun multireservoir system, Iran." Proceedings of the World Environ-838 mental & Water Resources Congress (ASCE EWRI 2014), Portland (Oregon). 839
- Brill., E., Flach, J., Hopkins, L., and Ranjithan, S. (1990). "MGA: A Decision Support 840 System for Complex, Incompletely Defined Problems." *IEEE Transactions on Systems*, 841 Man, and Cybernetics, 20(4), 745–757.
- Buhmann, M. (2003). Radial basis functions: theory and implementations. Cambridge uni-843 versity press Cambridge. 844

842

- Busoniu, L., Babuska, R., De Schutter, B., and Ernst, D. (2010). Reinforcement Learning 845 and Dynamic Programming Using Function Approximators. CRC Press, New York. 846
- Busoniu, L., Ernst, D., De Schutter, B., and Babuska, R. (2011). "Cross-Entropy Optimiza-847
- tion of Control Policies With Adaptive Basis Functions." IEEE Transactions on systems, 848 man and cybernetics-Part B: cybernetics, 41(1), 196–209. 849
- Castelletti, A., de Rigo, D., Rizzoli, A., Soncini-Sessa, R., and Weber, E. (2007). "Neuro-850 dynamic programming for designing water reservoir network management policies." Con-851 trol Engineering Practice, 15(8), 1031–1038. 852
- Castelletti, A., Galelli, S., Restelli, M., and Soncini-Sessa, R. (2010). "Tree-based rein-853 forcement learning for optimal water reservoir operation." Water Resources Research, 854 46(W09507).855
- Castelletti, A., Pianosi, F., Quach, X., and Soncini-Sessa, R. (2012). "Assessing water reser-856 voirs management and development in Northern Vietnam." Hydrology and Earth System 857 Sciences, 16(1), 189–199. 858
- Castelletti, A., Pianosi, F., and Restelli, M. (2013). "A multiobjective reinforcement learning 859 approach to water resources systems operation: Pareto frontier approximation in a single 860 run." Water Resources Research, 49. 861
- Castelletti, A., Pianosi, F., and Soncini-Sessa, R. (2008). "Water reservoir control under 862 economic, social and environmental constraints." Automatica, 44(6), 1595–1607. 863
- Castelletti, A., Pianosi, F., and Soncini-Sessa, R. (2012). "Stochastic and robust control 864

- of water resource systems: Concepts, methods and applications." System Identification,
   Environmental Modelling, and Control System Design, Springer, 383–401.
- <sup>867</sup> Celeste, A. and Billib, M. (2009). "Evaluation of stochastic reservoir operation optimization
  <sup>868</sup> models." Advances in Water Resources, 32(9), 1429–1443.
- <sup>869</sup> Chankong, V. and Haimes, Y. (1983). Multiobjective decision making: theory and methodol<sup>870</sup> oqy. North-Holland, New York, NY.
- <sup>871</sup> Chen, T. and Chen, H. (1995). "Universal approximation to nonlinear operators by neural
  <sup>872</sup> networks with arbitrary activation functions and its application to dynamical systems."

IEEE Transactions on Neural Networks, 6(4), 911–917.

- <sup>874</sup> Clark, E. (1950). "New York control curves." Journal of the American Water Works Association, 42(9), 823–827.
- Clark, E. (1956). "Impounding reservoirs." Journal of the American Water Works Association, 48(4), 349–354.
- <sup>878</sup> Coello Coello, C., Lamont, G., and Veldhuizen, D. V. (2007). Evolutionary Algorithms for
   <sup>879</sup> Solving Multi-Objective Problems (Genetic Algorithms and Evolutionary Computation).
   <sup>880</sup> Springer, New York, 2 edition.
- <sup>881</sup> Cohon, J. L. and Marks, D. (1975). "A review and evaluation of multiobjective programing
  techniques." Water Resources Research, 11(2), 208–220.
- <sup>883</sup> Cui, L. and Kuczera, G. (2005). "Optimizing water supply headworks operating rules un<sup>884</sup> der stochastic inputs: Assessment of genetic algorithm performance." Water Resources
  <sup>885</sup> Research, 41.
- <sup>886</sup> Cybenko, G. (1989). "Approximation by superpositions of a sigmoidal function." *Mathemat-*<sup>887</sup> *ics of control, signals and systems*, 2(4), 303–314.
- Dariane, A. and Momtahen, S. (2009). "Optimization of Multireservoir Systems Operation
  Using Modified Direct Search Genetic Algorithm." Journal of Water Resources Planning
  and Management, 135(3), 141–148.
- de Rigo, D., Castelletti, A., Rizzoli, A., Soncini-Sessa, R., and Weber, E. (2005). "A selec-

- tive improvement technique for fastening neuro-dynamic programming in water resources network management." *Proceedings of the 16th IFAC World Congress*, Prague (Czech Republic).
- <sup>895</sup> Deb, K. (2001). Multi-objective optimization using evolutionary algorithms. John Wiley &
  <sup>896</sup> Sons.
- <sup>897</sup> Deisenroth, M., Neumann, G., and Peters, J. (2011). "A Survey on Policy Search for <sup>898</sup> Robotics." *Foundations and Trends in Robotics*, Vol. 2, 1–142.
- <sup>899</sup> Desreumaux, Q., Côté, P., and Leconte, R. (2014). "Role of hydrologic information in
  <sup>900</sup> stochastic dynamic programming: a case study of the Kemano hydropower system in
  <sup>901</sup> British Columbia." *Canadian Journal of Civil Engineering*, 41(9), 839–844.
- Draper, A. and Lund, J. (2004). "Optimal Hedging and Carryover Storage Value." Journal
  of Water Resources Planning and Management, 130(1), 83–87.
- Esogbue, A. (1989). "Dynamic programming and water resources: Origins and interconnections." Dynamic Programming for Optimal Water Resources Systems Analysis, PrenticeHall, Englewood Cliffs.
- Faber, B. and Stedinger, J. (2001). "Reservoir optimization using sampling SDP with ensemble streamflow prediction (ESP) forecasts." *Journal of Hydrology*, 249(1), 113–133.
- Fleming, P., Purshouse, R., and Lygoe, R. (2005). "Many-Objective optimization: an engineering design perspective." *Proceedings of the Third international conference on Evolutionary Multi-Criterion Optimization*, Guanajuato, Mexico. 14–32.
- <sup>912</sup> Funahashi, K. (1989). "On the approximate realization of continuous mappings by neural
  <sup>913</sup> networks." Neural networks, 2(3), 183–192.
- <sup>914</sup> Gaggero, M., Gnecco, G., and Sanguineti, M. (2014). "Suboptimal Policies for Stochastic N-
- <sup>915</sup> Stage Optimization: Accuracy Analysis and a Case Study from Optimal Consumption."
- <sup>916</sup> Models and Methods in Economics and Management Science, F. E. Ouardighi and K.
- <sup>917</sup> Kogan, eds., number 198 in International Series in Operations Research & Management
- Science, Springer International Publishing, 27–50.

- Gass, S. and Saaty, T. (1955). "Parametric objective function Part II." Operations Research,
  3, 316–319.
- Giuliani, M., Galelli, S., and Soncini-Sessa, R. (2014). "A dimensionality reduction approach
  for Many-Objective Markov Decision Processes: application to a water reservoir operation
  problem." *Environmental Modeling & Software*, 57, 101–114.
- Giuliani, M., Herman, J., Castelletti, A., and Reed, P. (2014). "Many-objective reservoir
  policy identification and refinement to reduce policy inertia and myopia in water management." Water Resources Research, 50, 3355–3377.
- Giuliani, M., Mason, E., Castelletti, A., Pianosi, F., and Soncini-Sessa, R. (2014). "Universal approximators for direct policy search in multi-purpose water reservoir management: A
  comparative analysis." *Proceedings of the 19th IFAC World Congress*, Cape Town (South Africa).
- Gleick, P. and Palaniappan, M. (2010). "Peak water limits to freshwater withdrawal and
  use." *Proceedings of the National Academy of Sciences of the United States of America*,
  107(25), 11155–11162.
- <sup>934</sup> Guariso, G., Rinaldi, S., and Soncini-Sessa, R. (1986). "The Management of Lake Como: A
  <sup>935</sup> Multiobjective Analysis." Water Resources Research, 22(2), 109–120.
- <sup>936</sup> Guo, X., Hu, T., Zeng, X., and Li, X. (2013). "Extension of parametric rule with the hedging
  <sup>937</sup> rule for managing multireservoir system during droughts." *Journal of Water Resources*<sup>938</sup> *Planning and Management*, 139(2), 139–148.
- Hadka, D. and Reed, P. (2012). "Diagnostic assessment of search controls and failure modes
  in many-objective evolutionary optimization." *Evolutionary Computation*, 20(3), 423–452.
- Hadka, D. and Reed, P. (2013). "Borg: An Auto-Adaptive Many-Objective Evolutionary
  Computing Framework." *Evolutionary Computation*, 21(2), 231–259.
- Haimes, Y., Lasdon, L., and Wismer, D. (1971). "On a bicriterion formulation of the prob-
- lems of integrated system identification and system optimization." *IEEE Transactions on*
- <sup>945</sup> Systems, Man and Cybernetics, 1, 296–297.

- Hall, W. and Buras, N. (1961). "The dynamic programming approach to water-resources
  development." Journal of Geophysical Research, 66(2), 517–520.
- Heidrich-Meisner, V. and Igel, C. (2008). "Variable metric reinforcement learning methods
  applied to the noisy mountain car problem." *Recent Advances in Reinforcement Learning*,
  Springer, 136–150.
- <sup>951</sup> Hejazi, M. and Cai, X. (2009). "Input variable selection for water resources systems using a
  <sup>952</sup> modified minimum redundancy maximum relevance (mMRMR) algorithm." Advances in
  <sup>953</sup> Water Resources, 32(4), 582–593.
- <sup>954</sup> Herman, J. D., Reed, P. M., Zeff, H. B., and Characklis, G. W. (2015). "How Should
  <sup>955</sup> Robustness Be Defined for Water Systems Planning under Change?." Journal of Water
  <sup>956</sup> Resources Planning and Management.
- <sup>957</sup> Hornik, K., Stinchcombe, M., and White, H. (1989). "Multilayer feedforward networks are
  <sup>958</sup> universal approximators." *Neural networks*, 2(5), 359–366.
- <sup>959</sup> Kasprzyk, J., Reed, P., Kirsch, B., and Characklis, G. (2009). "Managing population and
  <sup>960</sup> drought risks using many-objective water portfolio planning under uncertainty." Water
  <sup>961</sup> Resources Research, 45(12).
- <sup>962</sup> Knowles, J. and Corne, D. (2002). "On metrics for comparing non-dominated sets." *Proceed-*<sup>963</sup> ings of the 2002 World Congress on Computational Intelligence (WCCI). IEEE Computer
  <sup>964</sup> Society, 711–716.
- <sup>965</sup> Koutsoyiannis, D. and Economou, A. (2003). "Evaluation of the parameterization<sup>966</sup> simulation-optimization approach for the control of reservoir systems." Water Resources
  <sup>967</sup> Research, 39(6), 1170–1187.
- Labadie, J. (2004). "Optimal operation of multireservoir systems: State-of-the-art review."
   Journal of Water Resources Planning and Management, 130(2), 93–111.
- 970 Loucks, D. and Sigvaldason, O. (1982). "Multiple-reservoir operation in North America.."
- <sup>971</sup> The operation of multiple reservoir systems, Z. Kaczmarck and J. Kindler, eds., IIASA
- 972 Collab. Proc. Ser., 1–103.

- <sup>973</sup> Loucks, D., van Beek, E., Stedinger, J., Dijkman, J., and Villars, M. (2005). Water Re<sup>974</sup> sources Systems Planning and Management: An Introduction to Methods, Models and
  <sup>975</sup> Applications. UNESCO, Paris, France.
- <sup>976</sup> Lund, J. and Guzman, J. (1999). "Derived operating rules for reservoirs in series or in
  <sup>977</sup> parallel." Journal of Water Resources Planning and Management, 125(3), 143–153.
- Maass, A., Hufschmidt, M., Dorfman, R., Thomas Jr, H., Marglin, S., and Fair, G. (1962).
  Design of water-resource systems. Harvard University Press Cambridge, Mass.
- Maier, H. and Dandy, G. (2000). "Neural networks for the prediction and forecasting of
  water resources variables: a review of modelling issues and applications." *Environmental*modelling & software, 15(1), 101–124.
- Maier, H., Kapelan, Z., Kasprzyk, J., Kollat, J., Matott, L., Cunha, M., Dandy, G., Gibbs,
- M., Keedwell, E., Marchi, A., Ostfeld, A., Savic, D., Solomatine, D., Vrugt, J., Zecchin,
- A., Minsker, B., Barbour, E., Kuczera, G., Pasha, F., Castelletti, A., Giuliani, M., and
- Reed, P. (2014). "Evolutionary algorithms and other metaheuristics in water resources:
- <sup>987</sup> Current status, research challenges and future directions ." *Environmental Modelling & Software*, 62(0), 271–299.
- Marbach, P. and Tsitsiklis, J. (2001). "Simulation-based optimization of Markov reward
   processes." *IEEE Transactions on Automatic Control*, 46(2), 191–209.
- Mayne, D. Q., Rawlings, J. B., Rao, C. V., and Scokaert, P. O. (2000). "Constrained model
  predictive control: Stability and optimality." *Automatica*, 36(6), 789–814.
- McDonald, R. I., Green, P., Balk, D., Fekete, B. M., Revenga, C., Todd, M., and Montgomery, M. (2011). "Urban growth, climate change, and freshwater availability." *Proceed- ings of the National Academy of Sciences*, 108(15), 6312–6317.
- <sup>996</sup> Mhaskar, H. and Micchelli, C. (1992). "Approximation by superposition of sigmoidal and <sup>997</sup> radial basis functions." *Advances in Applied mathematics*, 13(3), 350–373.
- <sup>998</sup> Momtahen, S. and Dariane, A. (2007). "Direct search approaches using genetic algorithms for <sup>999</sup> optimization of water reservoir operating policies." *Journal of Water Resources Planning*

- and Management, 133(3), 202-209.
- <sup>1001</sup> Moriarty, D., Schultz, A., and Grefenstette, J. (1999). "Evolutionary Algorithms for Rein-<sup>1002</sup> forcement Learning." *Journal of Artificial Intelligence Research*, 11, 199–229.
- Nalbantis, I. and Koutsoyiannis, D. (1997). "A parametric rule for planning and management
  of multiple-reservoir systems." Water Resources Research, 33(9), 2165–2177.
- <sup>1005</sup> Nicklow, J., Reed, P., Savic, D., Dessalegne, T., Harrell, L., Chan-Hilton, A., Karamouz,
- M., Minsker, B., Ostfeld, A., Singh, A., and Zechman, E. (2010). "State of the Art for
   Genetic Algorithms and Beyond in Water Resources Planning and Management." *Journal*
- <sup>1008</sup> of Water Resources Planning and Management, 136(4), 412–432.
- Oliveira, R. and Loucks, D. P. (1997). "Operating rules for multi reservoir systems." Water
   *Resources Research*, 33, 839–852.
- Park, J. and Sandberg, I. (1991). "Universal approximation using radial-basis-function networks." Neural computation, 3(2), 246–257.
- Peters, J. and Schaal, S. (2008). "Reinforcement learning of motor skills with policy gradients." *Neural networks*, 21(4), 682–697.
- <sup>1015</sup> Pianosi, F., Quach, X., and Soncini-Sessa, R. (2011). "Artificial Neural Networks and Multi
- Objective Genetic Algorithms for water resources management: an application to the Hoa
  Binh reservoir in Vietnam." *Proceedings of the 18th IFAC World Congress*, Milan, Italy.
- Piccardi, C. and Soncini-Sessa, R. (1991). "Stochastic dynamic programming for reservoir
   optimal control: Dense discretization and inflow correlation assumption made possible by
   parallel computing." Water Resources Research, 27(5), 729–741.
- Powell, W. (2007). Approximate Dynamic Programming: Solving the curses of dimensional *ity.* Wiley, NJ.
- Reed, P., Hadka, D., Herman, J., Kasprzyk, J., and Kollat, J. (2013). "Evolutionary Multiobjective Optimization in Water Resources: The Past, Present, and Future." Advances in
  Water Resources, 51, 438–456.
- <sup>1026</sup> Reed, P. M. and Kollat, J. B. (2013). "Visual analytics clarify the scalability and effective-

- ness of massively parallel many-objective optimization: A groundwater monitoring design
  example." Advances in Water Resources, 56, 1–13.
- ReVelle, C. and McGarity, A. E. (1997). Design and operation of civil and environmental
   engineering systems. John Wiley & Sons.
- <sup>1031</sup> Rippl, W. (1883). "The capacity of storage reservoirs for water supply." *Minutes of the*
- <sup>1032</sup> Proceedings, Institution of Civil Engineers, Vol. 71, Thomas Telford. 270–278.
- Rosenstein, M. and Barto, A. (2001). "Robot weightlifting by direct policy search." Inter national Joint Conference on Artificial Intelligence, Vol. 17, Citeseer. 839–846.
- <sup>1035</sup> Sehnke, F., Osendorfer, C., Rückstieß, T., Graves, A., Peters, J., and Schmidhuber, J. (2010).
- <sup>1036</sup> "Parameter-exploring policy gradients." *Neural Networks*, 23(4), 551–559.
- <sup>1037</sup> Soncini-Sessa, R., Castelletti, A., and Weber, E. (2007). Integrated and participatory water
   <sup>1038</sup> resources management: Theory. Elsevier, Amsterdam, NL.
- Sutton, R., McAllester, D., Singh, S., and Mansour, Y. (2000). "Policy Gradient Methods for
   Reinforcement Learning with Function Approximation." Advances in Neural Information
   Processing Systems, 12, 1057–1063.
- Tejada-Guibert, J., Johnson, S., and Stedinger, J. (1995). "The value of hydrologic information in stochastic dynamic programming models of a multireservoir system." Water *Resources Research*, 31(10), 2571–2579.
- Tikk, D., Kóczy, L., and Gedeon, T. (2003). "A survey on universal approximation and its
  limits in soft computing techniques." *International Journal of Approximate Reasoning*,
  33(2), 185–202.
- Tsitsiklis, J. and Van Roy, B. (1996). "Feature-Based Methods for Large Scale Dynamic
  Programming." *Machine Learning*, 22, 59–94.
- Tu, M., Hsu, N., and Yeh, W. (2003). "Optimization of reservoir management and operation
  with hedging rules." Journal of Water Resources Planning and Management, 129(2), 86–
  97.
- 1053 U.S. Army Corps of Engineers (1977). Reservoir System Analysis for Conservation, Hy-

- drologic Engineering Methods for Water Resources Development. Hydrologic Engineering
   Center, Davis, CA.
- Whiteson, S. and Stone, P. (2006). "Evolutionary function approximation for reinforcement
  learning." The Journal of Machine Learning Research, 7, 877–917.
- Whitley, D., Dominic, S., Das, R., and Anderson, C. (1994). Genetic reinforcement learning
   for neurocontrol problems. Springer.
- Woodruff, M., Reed, P., and Simpson, T. (2013). "Many objective visual analytics: rethinking the design of complex engineered systems." *Structural and Multidisciplinary Optimiza- tion*, 1–19.
- Yeh, W. (1985). "Reservoir management and operations models: a state of the art review."
   Water Resources Research, 21 (12), 1797–1818.
- Zitzler, E., Thiele, L., Laumanns, M., Fonseca, C., and da Fonseca, V. (2003). "Performance
   assessment of multiobjective optimizers: an analysis and review." *IEEE Transactions on Evolutionary Computation*, 7(2), 117–132.
- <sup>1068</sup> Zoppoli, R., Sanguineti, M., and Parisini, T. (2002). "Approximating networks and extended
- ritz method for the solution of functional optimization problems." Journal of Optimization
- 1070 Theory and Applications, 112(2), 403-440.

# 1071 List of Figures

1	Schematization of the evolutionary multi-objective direct policy search (EMODPS	5)
	approach. The dashed line represents the model of the system, the gray box	
	the MOEA algorithm.	44
2	(a) Map of the Red River basin and (b) schematic representation of the main	
	components described in the model.	45
3	Policy performance obtained with different structures of ANNs and RBFs	
	over the optimization horizon 1962-1969 (a), and evaluation of the associated	
	Pareto fronts in terms of generational distance (b), additive $\varepsilon$ -indicator (c),	
	and hypervolume indicator (d).	46
4	Probability of attainment with a threshold equal to $75\%$ (a) and to $95\%$ (b)	
	of the best metric values for different ANN and RBF architectures	47
5	Analysis of runtime search dynamics for ANN and RBF operating policy op-	
	timization in terms of generational distance (a), additive $\varepsilon$ -indicator (b), and	
	hypervolume (c). $\ldots$	48
6	Validation of EMODPS operating policies via comparison of ANN and RBF	
	performance over the optimization and the validation horizons (a) and com-	
	parison with SDP solutions (b)	49
7	Visualization of the compromise operating policy selected in Figure 5b	50
	1 $     2 $ $     3 $ $     4 $ $     5 $ $     6 $ $     7$	<ol> <li>Schematization of the evolutionary multi-objective direct policy search (EMODPS approach. The dashed line represents the model of the system, the gray box the MOEA algorithm.</li> <li>(a) Map of the Red River basin and (b) schematic representation of the main components described in the model.</li> <li>Policy performance obtained with different structures of ANNs and RBFs over the optimization horizon 1962-1969 (a), and evaluation of the associated Pareto fronts in terms of generational distance (b), additive ε-indicator (c), and hypervolume indicator (d).</li> <li>Probability of attainment with a threshold equal to 75% (a) and to 95% (b) of the best metric values for different ANN and RBF operating policy optimization in terms of generational distance (a), additive ε-indicator (b), and hypervolume (c).</li> <li>Analysis of runtime search dynamics for ANN and RBF operating policy optimization in terms of generational distance (a), additive ε-indicator (b), and hypervolume (c).</li> <li>Validation of EMODPS operating policies via comparison of ANN and RBF performance over the optimization and the validation horizons (a) and comparison with SDP solutions (b).</li> <li>Visualization of the compromise operating policy selected in Figure 5b.</li> </ol>

FIG. 1. Schematization of the evolutionary multi-objective direct policy search (EMODPS) approach. The dashed line represents the model of the system, the gray box the MOEA algorithm.



FIG. 2. (a) Map of the Red River basin and (b) schematic representation of the main components described in the model.



FIG. 3. Policy performance obtained with different structures of ANNs and RBFs over the optimization horizon 1962-1969 (a), and evaluation of the associated Pareto fronts in terms of generational distance (b), additive  $\varepsilon$ -indicator (c), and hypervolume indicator (d).



(a) Policy Performance with different ANNs and RBFs architectures



FIG. 4. Probability of attainment with a threshold equal to 75% (a) and to 95% (b) of the best metric values for different ANN and RBF architectures.

FIG. 5. Analysis of runtime search dynamics for ANN and RBF operating policy optimization in terms of generational distance (a), additive  $\varepsilon$ -indicator (b), and hypervolume (c).



# FIG. 6. Validation of EMODPS operating policies via comparison of ANN and RBF performance over the optimization and the validation horizons (a) and comparison with SDP solutions (b).











# FIG. 7. Visualization of the compromise operating policy selected in Figure 6b.