

Policy Representation Learning for multiobjective reservoir policy design with different objective dynamics

Marta Zaniolo¹, Matteo Giuliani¹, and Andrea Castelletti¹

¹Department of Electronics, Information and Bioengineering, Politecnico di Milano, Milan, Italy

Key Points:

- We introduce a novel method to define an optimal input set for a multipurpose dam operating policy that varies with the objective trade-off.
- Better informed policies are able to mitigate conflicts between water users and achieve system-wide benefits.
- The addition of information in policy design increases the policies robustness towards extreme hydrological conditions.

Corresponding author: Marta Zaniolo, marta.zaniolo@polimi.it

This article has been accepted for publication and undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1029/2020WR029329](https://doi.org/10.1029/2020WR029329).

This article is protected by copyright. All rights reserved.

Abstract

Most water reservoir operators make use of forecasts to inform their decisions and enhance water systems flexibility and resilience by anticipating hydrological extremes. Yet, despite numerous candidate hydro-meteorological variables and forecast horizons may potentially be beneficial to operations, the best information set for a given problem is often not evident. Additionally, in multi-purpose systems characterized by multiple demands with varying vulnerabilities and temporal scales, this information set might change according to the objective tradeoff. In this work, we contribute a novel method to learn the optimal policy representation (i.e., policy input set) by combining a feature selection routine with a multi-objective Direct Policy Search framework in order to retrieve the best policy input set online (i.e., while learning the policy) and dynamically with the objective trade-off. The selected policy search routine is the Neuro-Evolutionary Multi-Objective Direct Policy Search (NEMODPS) which generates flexible policy shapes adaptive to online changes in the input set. This approach is demonstrated on the case study of Lake Como (Italy), where the operating objectives are highly heterogeneous in their dynamics (fast and slow) and vulnerabilities (wet and dry extremes). We show how varying objectives, and tradeoffs therein, benefit from a different policy representation, ultimately yielding remarkable results in terms of conflict mitigation between different users. More informed policies, moreover, show higher robustness when re-evaluated across a suite of different hydrological conditions.

1 Introduction

Water reservoirs have long been fundamental components of coupled human-water systems worldwide, providing communities with green and affordable electricity, water supply for agricultural and urban consumption, and flood protection. Yet, lately, new concerns are arising regarding the reliability of water systems as climate change increases the likelihood of extreme events, and economic development exacerbates water demands and conflicts (Fletcher *et al.*, 2019; Herman *et al.*, 2020). One way of increasing resilience and reliability of water systems is to build more, larger, infrastructures, however, this hard path to capacity expansion is costly and often yields unintended cross-sectoral externalities (Gleick, 2003). An alternative, soft-path towards resilience advocates the improvement of the operating rules used to control the existing water infrastructures to enhance their capability to anticipate weather extremes, and timely prepare for them.

Traditionally, the operating policy of water reservoirs was conditioned upon very limited information systems comprising reservoir storage and a cyclostationary time index (Hejazi *et al.*, 2008). More recently, Turner *et al.* (2019) showed that most water system operators across the US make use of streamflow forecasts to further improve operations. The employed forecast horizon is however reservoir-specific, and, when official guidelines are absent, operators seem to rely on their expert judgment to identify their forecast horizon of choice. In the water resources literature, few studies have tackled the issue of the optimal selection of streamflow forecast horizon for a single-objective reservoir operated for water supply (Anghileri *et al.*, 2016), hydropower (Hamlet *et al.*, 2002; Block, 2011; Xu *et al.*, 2014), or for a generic concave objective function (Zhao *et al.*, 2014, Zhao *et al.*, 2019). Additionally, the breath of information sources that was demonstrated to be valuable to inform reservoir operations is by no means limited to streamflow forecasts, but includes the previous period's inflow (Gal, 1979; Maidment and Chow, 1981), available hydrological observations (Denaro *et al.*, 2017), traditional (Hejazi and Cai, 2011) or basin-specific (Zaniolo *et al.* 2018, Zaniolo *et al.*, 2019) drought indexes, measures of snow abundance (Desreumaux *et al.*, 2014; Giuliani *et al.*, 2016a), shifts in hydrological regimes (Turner and Galelli, 2016), teleconnection indices (Libisch-Lehner *et al.*, 2019), and sea surface temperature measured in appropriate locations (Giuliani *et al.*, 2019; Zaniolo *et al.*, 2021a).

64 While these studies are a great demonstration of the potential of using unconven-
65 tional policy representations in policy design, none of them attempts at automatizing
66 representation learning in a portable framework. Additionally, no attention has been given
67 to a major challenge to learning an optimal policy representation, i.e., the coexistence
68 of multiple operating objectives. In fact, previous studies either consider systems oper-
69 ated for a single purpose (i.e., reservoirs operated just for hydropower), or specify only
70 one policy representation for the entire tradeoff space. In multi-purpose water reservoir
71 systems, however, common operating targets, e.g., flood protection and water supply, can
72 be vastly heterogeneous in their dynamics and vulnerabilities. Flood events are gener-
73 ally caused by the onset of fast and intense wet meteorological extreme events, while wa-
74 ter supply failures are the result of a prolonged period of water shortage caused by slow-
75 developing dry hydrological extremes, i.e., droughts. In these systems, defining an ap-
76 propriate policy representation becomes more intricate. On the one hand, a flood-conservative
77 policy benefits from a short lead time look-ahead information that conveys peak inflow
78 magnitude and timing, on the other, a water supply-prone policy seeks predictors that
79 are relevant for the onset of a prolonged water shortage to timely activate hedging strate-
80 gies. The tradeoff space between these two opposite solutions is populated by an ensem-
81 ble of policies diversely balancing opposite control targets. Such behavior is shown for
82 a fixed policy representation via sensitivity analysis to policy inputs for alternative trade-
83 offs (Quinn *et al.*, 2017; Doering *et al.*, 2021).

84 In this work, for the first time, we hypothesize and quantitatively demonstrate that
85 in Multi-Objective (MO) problems different objective tradeoffs require different infor-
86 mation, and selection of policy representation should be tradeoff-specific. Our results demon-
87 strate that one policy input set is inadequate to represent the entire space of different
88 control behaviors that may emerge for alternative tradeoffs.

89 Part of the reason why a tradeoff dynamic selection was never performed is that
90 traditional policy search routines only support static and prespecified input sets, thereby
91 not allowing the evolution of a population with heterogeneous input sets. In this work,
92 we propose a novel technique to automatically learn a Pareto front of optimal policies
93 and their representations for a multipurpose water system. The method is applicable to
94 large and heterogeneous datasets of candidate policy inputs, from meteorological and hy-
95 drological forecasts with disparate horizons, to observational data. The framework, namely
96 SINEPS, Selection of Information for NeuroEvolutionary Policy Search, combines au-
97 tomatic feature selection with NEMODPS (NeuroEvolutionary Multi-Objective Direct
98 Policy Search, Zaniolo *et al.* (2021b)), a policy search routine that can accomodate changes
99 in the policy input set. SINEPS starts with a simple operating policy and a minimal pol-
100 icy representation and gradually includes new inputs to the policy representations while
101 automatically adjusting the policy processing capacity. For every Pareto efficient pol-
102 icy, the selected input is the one that explains most of the information gap between the
103 policy itself, and an ideal, deterministic, Perfect Operating Policy, designed under the
104 assumption of perfect knowledge of future disturbance.

105 This framework is tested on the real-world case study of the multi-purpose Lake
106 Como, operated to meet two conflicting and heterogeneous objectives of flood protec-
107 tion and water supply, mainly for irrigation. The flood objective is characterized by fast
108 dynamics and vulnerability towards wet extremes, while irrigation supply is character-
109 ized by a slow dynamic and vulnerability towards dry extremes. In this paper, the dataset
110 of candidate policy inputs is composed of perfect streamflow forecasts at different lead
111 times.

112 1.1 Literature review on policy representation learning

113 The problem of learning a policy representation is not unique to water resources
114 management, on the contrary, it is widely addressed in the control community, finding

115 applications in diverse fields, from spatial path scheduling (Whiteson *et al.*, 2005), stock
116 index trading (Si *et al.*, 2017), to virtually any autonomous robot control task (e.g., Hachiya
117 and Sugiyama, 2010; Lesort *et al.*, 2018). In this section, we propose a literature review
118 on policy representation learning that goes beyond the existing experience in dam pol-
119 icy design in order to present and discuss the wider background and challenges the in-
120 spired the design of SINEPS, and motivate its algorithmic choices.

121 When designing an operating policy for a given system, defining the policy repre-
122 sentation corresponds to selecting its input set. Such problem is generally tackled by pair-
123 ing Feature Extraction with Policy Search (Liu *et al.*, 2015; Lesort *et al.*, 2018). Feature
124 Extraction refer to a family of techniques that transform an original dataset into a more
125 compact, while still highly informative dataset (Cunningham, 2008). Policy Search meth-
126 ods aim at learning an optimal operating policy for a system (e.g., a release policy from
127 a reservoir) with respect to its objective functions (e.g., flood and water supply). In the
128 proposed taxonomy, we identify *a priori*, *a posteriori*, and *online* approaches to pairing
129 feature extraction and policy search for learning a policy and its representation.

130 In the first *a priori* approach, the feature extraction step is antecedent and inde-
131 pendent from the policy search step. First, the feature extraction routine reduces the
132 dimensionality of the dataset of candidate features for example extracting few relevant
133 features from the dataset, removing irrelevant ones, or generating new features by ap-
134 propriately combining existing ones. The reduced dataset represents the selected pol-
135 icy representation, and is used for policy search. The dimension reduction is generally
136 achieved via i) data compression techniques, e.g., autoencoders (e.g., Morimoto *et al.*,
137 2008), or Principal Component Analysis (Nouri and Littman, 2010), that map the ini-
138 tial dataset into a lower dimensional latent space that retains most of its information con-
139 tent, ii) using a target control sequence to identify relevant policy drivers (Kroon and
140 Whiteson, 2009; Giuliani *et al.*, 2015; Denaro *et al.*, 2017), or, iii) via expert-based fea-
141 ture selection (e.g., Akrouf *et al.*, 2012) or extraction (e.g., Sturtevant and White, 2006;
142 Giuliani and Castelletti, 2019) to design a problem-specific representation. In general,
143 a priori approach to policy representation is advisable whenever there is sufficient knowl-
144 edge of the task to confidently devise an appropriate feature set. This very low compu-
145 tationally demanding approach, in fact, does not offer any guarantees on the optimal-
146 ity of the chosen representation (Lesort *et al.*, 2018).

147 The *a posteriori* approach evaluates the suitability of a policy representation by
148 assessing the performance of the policy conditioned upon it. Multiple policies are designed
149 with alternative input sets, and the desired representation is identified as the one gen-
150 erating the best performing policy. In principle, the entire combinatorial space of fea-
151 tures subsets could be exhaustively explored, yielding to an optimal solution albeit re-
152 sulting computationally intractable for non-trivial datasets (see, e.g., Gaudel and Sebag,
153 2010). Alternatively, for modest datasets, hill-climbing approaches incrementally add fea-
154 tures to the representation retaining the most successful ones (Wright *et al.*, 2012; Zhang,
155 2009; Tan *et al.*, 2013). Finally, an initial *a priori* reduction can be applied to select a
156 limited number of candidate representations that are then exhaustively compared *a pos-*
157 *teriori* (Giuliani *et al.*, 2016a; Castelletti *et al.*, 2016). In general, *a posteriori* feature
158 representation is significantly more computationally burdensome than the *a priori* coun-
159 terpart. Yet, an exhaustive *a posteriori* search can be performed with virtually no pre-
160 existing knowledge of the task, and guarantees the optimality of the derived feature rep-
161 resentation. Both *a priori*, and *a posteriori* approaches in general rely on heavy expert-
162 based manual engineering in defining potentially appropriate policies representations to
163 implement or test (Bengio *et al.*, 2013).

164 The third, *online* approach, interleaves feature extraction phases throughout the
165 policy search process, using progressively refined feature representations to support pol-
166 icy learning. Representations are updated during the search via supervised learning, by
167 extracting features that approximate the state space (Curran *et al.*, 2016; Alvernaz and

168 *Togelius*, 2017), state-transition space (*Assael et al.*, 2015; *Van Hoof et al.*, 2016), or the
 169 reward trajectory (*Munk et al.*, 2016; *Oh et al.*, 2017) of the policy learned thus far (for
 170 a comprehensive review, see *Lesort et al.*, 2018). The adjusted representation is then em-
 171 ployed to refine policy search in a feedback loop between the two routines. Computa-
 172 tionally, *online* approaches are more expensive than *a priori*, but less than *a posteriori*
 173 methods, while handling significantly larger datasets of candidate information.

174 1.2 SINEPS

175 In this work, we present a novel method for *online* dynamic policy representation
 176 called SINEPS, Selection of Information for NeuroEvolutionary Policy Search. It requires
 177 the selection of i) a feature extraction method, ii) a policy search routine, and iii) a strat-
 178 egy to interface the two.

- 179 1. **Feature extraction method:** Several *online* policy representation routines em-
 180 ploy Feature extraction techniques that reduce the dimensionality of the repre-
 181 sentation by projecting the initial feature space into a lower dimensional latent
 182 space that preserves information content. However, such an approach does not guar-
 183 antee that any candidate feature is actually excluded from the problem formula-
 184 tion (*Loscalzo et al.*, 2015). As a result, while the operating policy can actually
 185 benefit from a lower-dimensional representation, the actual problem size remains
 186 unchanged. In an operational setting, this implies that the entire dataset of ini-
 187 tial features must be retrieved continuously. Alternatively, Feature Selection meth-
 188 ods are a subset of the feature extraction techniques that reduces the dataset size
 189 by identifying a subset of the initial features. Some authors suggest the use of fea-
 190 ture selection routines, rather than information encoders, for representation learn-
 191 ing, in order to effectively restrict the number of candidate variables included in
 192 the problem formulation (e.g., *Loscalzo et al.*, 2015). The representation obtained
 193 through variable selection, moreover, highlights relevant policy drivers, is easily
 194 interpretable, and can thus generate insights on the task at hand. Within Feature
 195 Selection techniques, the iterative online framework can accommodate simple correla-
 196 tion-based variable filtering (i.e., the variables that are most correlated with the tar-
 197 get are selected), as well as non-linear model-based selection routines (e.g., IIS,
 198 *Castelletti et al.*, 2010). Here, we use a correlation-based filtering approach, where
 199 the correlation is measured in Symmetric Uncertainty (SU, *Blum and Langley*,
 200 1997). SU is a normalized version of the Mutual Information metric (MI, *Shan-*
 201 *non*, 1948) that quantifies the degree of similarity between two variables, or, more
 202 specifically, the amount of information that can be obtained on one variable by
 203 observing the other. Entropy-based techniques like SU are model-free and gen-
 204 eralizable to any modeling context, as they do not require to assume any functional
 205 relationship between the variables (*MacKay*, 2003), contrary to simpler metrics
 206 such as correlation coefficients that assume a linear dependence. The use of SU
 207 is supported in the information theoretic literature and was demonstrated to out-
 208 perform several other feature selection methods on a suite of 15 benchmark fea-
 209 ture selection problems (*Zhang and Chen*, 2021).

210 Note that SU is employed as a screening tool that allows to detect promising pol-
 211 icy representations by identifying candidate variables with high information con-
 212 tent across different objectives. This is intended to avoid an exhaustive approach
 213 that would test every possible candidate representation in policy search, which would
 214 be computationally untractable. The policy search step, described below, evolves
 215 policies with different representations to generate a Pareto front of optimal poli-
 216 cies and applies further selection pressure onto alternative policy representations
 217 thereby further refining the representation selection in a policy search context.

- 218 2. **Policy Search Method:** Direct Policy Search (DPS) is emerging as one of the
 219 most effective, and widely applied methods to design optimal operating policies

for multi-purpose reservoir operations, given its multi-objective nature, flexibility in problem and objective formulation, and data-driven nature that allows to use trajectories of non-modeled information in policy design (Giuliani *et al.*, 2016b). DPS defines the operating policy within a prespecified class of functions and solves a problem of optimal functional parameterization with respect to the problem's objectives (Zatarain *et al.*, 2017; Quinn *et al.*, 2018; Giuliani *et al.*, 2019; Quinn *et al.*, 2019). Flexible universal approximators such as Neural Networks (NNs) are generally employed to parameterize the operating policy in order not to restrict the parametrical search to a small functional subspace that may not contain skillful solutions (Giuliani *et al.*, 2014, Giuliani *et al.*, 2018). The architecture of a NN employed for policy design includes as many input nodes as the number of features in the policy representation, and as many output nodes as the decisions to be taken on the system, e.g., reservoir release decisions. Finally, the internal NN complexity, i.e., number of hidden nodes, connections, and layers, is crucial to determine the network processing capability and training requirements. The a priori definition of the optimal network complexity for a given problem would require a perfect knowledge of the operational task, which is in general unavailable. Therefore, in practical application, the network architecture is selected by the modeler via few manual trials and errors balancing the network approximation capacity, training costs, and overfitting tendency. Given its rigid, prespecified, policy structure, DPS techniques do not support dynamic changes in the dimensionality of the policy feature representation.

A promising alternative that obviates to policy rigidity is represented by NeuroEvolution (NE), a set of techniques that employs evolutionary algorithms to evolve neural networks in terms of their architectures and parameters. These techniques generally begin with a population of simple networks and progressively build more sophisticated ones by applying new architectural elements (nodes and connections). The evolutionary competition ultimately determines the optimal network complexity. By pairing NE with DPS, it is possible to derive policy search routines that support online changes in policy architecture. Popular NE algorithms (e.g., NEAT Stanley and Miikkulainen, 2002) are, however, strictly applicable to single-objectives problems. The here employed NeuroEvolutionary Multi-Objective Direct Policy Search (NEMODPS), is the first NE routine specifically designed to solve MO problems in one algorithmic iteration (Zaniolo *et al.*, 2021b). NEMODPS is here employed for the first time to jointly evolve policies with different feature representations. In general, not all the policy representations identified in the feature selection step will survive the evolution pressure, thus refining the selection of optimal representations via policy competition. NEMODPS will be briefly introduced in Section 2.2 of the Methods. The reader is referred to Zaniolo *et al.* (2021b) for a more detailed analysis of NEMODPS, and its benchmarking against traditional DPS in terms of performance and computational costs.

- 3. Interfacing strategy:** in many applications, the selection of relevant features is performed via supervised learning using as target the state, state-transition, state-value spaces, or the cost trajectory produced by the policy learned thus far (for a review, see Lesort *et al.*, 2018). Cost-based selection is generally recognized as more effective in identifying task-oriented policy representations (Loscalzo *et al.*, 2015), however, in multi-objective problems, the coexistence of multiple cost signals complicates the cost-based selection process. In SINEPS, we propose a novel interfacing strategy that is both task-tailored, and suitable for MO problems. In particular, we use as reference a deterministic Perfect Operating Policy (POP) that assumes full knowledge of future system disturbance. For a given state, we contrast the actions extracted from the POP to those extracted from the policy under design. We assume that the difference in actions is due to the information gap in the policies representations, and thus surrogates the information that the designed policy would require to meet the POP performance. The trajectory of ac-

275 tion residuals is used as an interfacing strategy, and employed as target for fea-
 276 ture selection. Such target can be considered task-relevant, as it is a proxy of the
 277 policy information deficiency for a given task. Additionally, it can be applied to
 278 MO problems by contrasting each Pareto efficient policy with the corresponding
 279 perfect counterpart supporting a tradeoff dynamic feature selection.

280 To summarize, SINEPS combines feature selection, neuroevolution and an original in-
 281 terfacing strategy. The choices made in the selection and development of the building
 282 tools of SINEPS target the overarching goal of designing the first multi-objective fea-
 283 ture representation learning routine that automatically specifies an optimal policy rep-
 284 resentations for each tradeoff.

285 This paper is organized as follows. The next section presents the methods of this work,
 286 by presenting the methodological Framework 2.1, and expanding on the key concepts and
 287 tools employed in the methodology, including NEMODPS 2.2. Section 3 is dedicated to
 288 the presentation of the case study and experimental settings. Results are discussed in
 289 Section 4, and in the following Section 5 we draw conclusions and introduce some dis-
 290 cussion points.

291 2 Methods

292 In this work, we consider a water reservoir system modeled as a discrete-time, pe-
 293 riodic, non-linear, stochastic process defined by a state variable s_t (reservoir storage),
 294 a control variable u_t representing the release decision from the dam gates, stochastic dis-
 295 turbances ε_{t+1} (net reservoir inflow), and a state-transition function $f(\cdot)$: $s_{t+1} = f(s_t, r_{t+1}, \varepsilon_{t+1})$
 296 where the effective release r_{t+1} coincides with the release decision u_t corrected, where
 297 appropriate, with a non-linear release function $R_t(s_t, u_t, \varepsilon_{t+1})$ determining the minimum
 298 and maximum releases feasible for the time interval $[t, t+1)$ to respect physical and le-
 299 gal constraints. The operating policy π determines the release decision from the water
 300 reservoir $u_t = \pi(\cdot)$ at each time step t over the simulation horizon H . The objective
 301 of this work is to design the optimal operating policy and relative representation for this
 302 system by solving a minimization problem formulated as follows:

$$\min_{\pi, \mathbf{I}_t, \zeta(\theta)} \mathbf{J}(\pi, s_0, \varepsilon_1^H) \quad (1)$$

303 where we search the minimum of the multidimensional objective function \mathbf{J} , here
 304 interpreted as cost, with respect to the closed loop operating policy π , its representa-
 305 tion \mathbf{I}_t , functional class ζ and relative parameterization θ . In particular, the operating
 306 policy π is conditioned upon basic information (i.e., the reservoir storage s_t a time in-
 307 dex d_t), and an additional vector of information \mathbf{I}_t searched within the dataset of candi-
 308 date information as in $\pi = \pi(s_t, d_t, \mathbf{I}_t)$. Among the available policy search methods,
 309 parametric approaches define π within a class of functions ζ , and search its optimal pa-
 310 rameterization θ . The employed NEMODPS technique supports the conjunct search of
 311 the optimal functional class ζ and relative parameters as in $\zeta(\theta)$.

312 In general, in MO problems, conflicts occur between different operating objectives,
 313 and the solution is constituted by a set of non-dominated (or Pareto optimal) solutions
 314 $\mathcal{P}^* = \{\pi^* | \nexists \pi \prec \pi^*\}$, which maps onto the Pareto front $\mathcal{F}^* = \{\mathbf{J}(\pi^*, \mathbf{x}_0, \varepsilon_1^H) | \pi^* \in$
 315 $\mathcal{P}^*\}$. For a more complete problem formulation please refer to the *Detailed Problem For-*
 316 *mulation* section of the Supplementary Information.

317 2.1 Framework

318 In this section, we present the flowchart of the proposed SINEPS framework em-
 319 ployed to approach Problem 1, reported in Figure 1 and organized in numbered blocks.

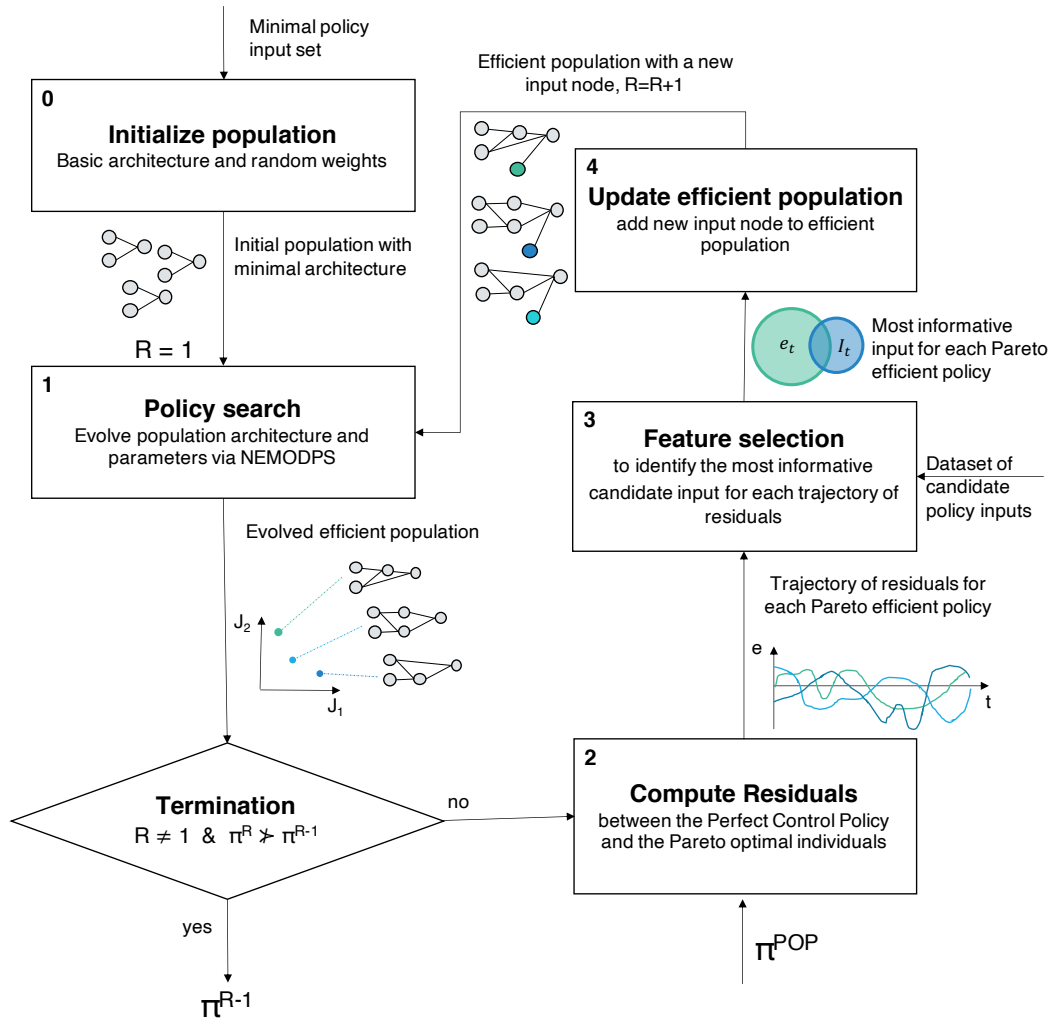


Figure 1. SINEPS flowchart. By looping through the building blocks of this flowchart, the procedure complexifies the initial population in terms of feature representation and policy architecture.

320 **0:** The procedure begins in round R1, with the initialization of a population of sim-
 321 ple neural networks, a minimal architecture, and random weights. At this stage, the pol-
 322 icy representation is also minimal, comprising a cyclostationary time index d_t and the
 323 reservoir storage s_t , namely, $\pi^{R1} = \pi^{R1}(d_t, s_t)$.

324 **1:** This population is the input to the Policy Search building block that employs
 325 NEMODPS. For a given input set, NEMODPS evolves policies' architecture and param-
 326 eters in a MO problem (more details in the dedicated Subsection 2.2). The output of this
 327 step is an ensemble of Pareto efficient operating policies, each specified with a tailored
 328 architecture, resulting in an architecturally heterogeneous population.

329 **2:** In the first round, the flowchart proceeds to the building block named Compute
 330 Residuals. In this step, we contrast the operating decisions produced by each Pareto ef-
 331 ficient policy with the decisions given by a Perfect Operating Policy (POP) extracting
 332 the trajectories of decision residuals e_t , i.e., the difference in the decisions selected by
 333 the minimally informed policy under design π^{R1} , and the perfectly informed policy π^{POP} .
 334 The calculated residuals are assumed to be due to their information gap (more details
 335 in the dedicated Section 2.3).

336 **3:** In the Feature Selection step, we search the dataset of candidate policy inputs
 337 D to identify the most informative feature for π^{R1} . For this purpose, we compute the
 338 SU metric between each vector of residual trajectory in e_t , i.e., the policy's information
 339 gap, and the candidate policy input dataset D . SU quantifies the amount of informa-
 340 tion shared between e_t and each candidate input, allowing to identify the most promis-
 341 ing feature representation by selecting the feature that explains most of the policy in-
 342 formation gap. SU is defined in [0,1] and can be computed for two variables X and Y
 343 as:

$$SU(X, Y) = 2 * \frac{MI(X, Y)}{H(X) + H(Y)} = 2 * \frac{H(X) + H(Y) - H(X, Y)}{H(X) + H(Y)} \quad (2)$$

344 Where $H(X)$ and $H(Y)$ are the entropy of the variable X and Y, and $H(X, Y)$ is their
 345 joint entropy.

346 Because the trajectory of residuals is computed independently for each efficient pol-
 347 icy, the inputs selected are policy-specific, and may vary across the tradeoff space.

348 **4:** Each efficient policy is then updated by including the selected feature in the in-
 349 put set, with a single input-output connection and a randomly initialized weight. The
 350 population of policies is now heterogeneous in its feature representation. Such popula-
 351 tion will now enter round R2 of SINEPS, with an update representation that includes
 352 the tailored information I_t , $\pi^{R2} = \pi^{R2}(d_t, s_t, I_t)$. In step 1 of the second round R2, this
 353 population is further evolved via NEMODPS. Individuals will appropriately complex-
 354 ify their architecture by genetic evolution to adapt to the newly inserted input, and learn
 355 how to make use of its information content. Neuro-evolutionary competition will further
 356 filter feature representation, causing only the fittest representations to survive in the ef-
 357 ficient policies of round R2. Note that this framework performs a joint optimization of
 358 the policy inputs and architecture which cannot be decoupled. In particular, the input
 359 layer contains the information that a policy can access, while the policy structure gov-
 360 erns how this information is used and translated into a control decision. Therefore, on
 361 the one hand, a policy with an inadequate input layer won't be able to make good con-
 362 trol decisions because poorly informed, no matter how well the policy structure can trans-
 363 late input into decisions. Similarly, when a new input is added to an existing policy, the
 364 structural optimization is necessary for the policy to learn how to use the input, i.e., to
 365 build the structural elements (connection, nodes), that will enable it to appropriately
 366 use it to make more informed control decisions. Without a structural optimization, the
 367 new input would be unused.

SINEPS proceeds analogously until the Termination check is positive, namely when the efficient Pareto set at Round R does not significantly dominate the Pareto set in the previous round: $\pi^{*R} \not\prec \pi^{*R-1}$. More details on the termination criterion are presented in Section 2.4. Upon termination, we retain as efficient solutions the Pareto set generated at the previous round $R-1$, as it achieves virtually the same performance as round R with a simpler representation.

2.2 NEMODPS

In this section, we give an overview of the main components of NEMODPS, the policy search routine employed in this study. NEMODPS builds on a recent Reinforcement Learning branch called Neuro-Evolution (NE) (*Stanley and Miikkulainen, 2003; Floreano et al., 2008*), which employs Evolutionary Algorithms to optimize neural network architectures and parameters. NEMODPS algorithm is inspired by NEAT (*Stanley and Miikkulainen, 2002*), and the subsequent literature of NEAT improvements targeting complex control problems, vast decision spaces, and noisy environments. Additionally, NEMODPS contains original strategies to address the specific complexities of multi-objective optimization problems, which make NEMODPS the first multi-objective NE algorithm. An in-depth explanation of NEMODPS can be found in Zaniolo, 2021, but here we discuss the main algorithmic components.

Key elements of NEMODPS are (1) a process of evolutionary complexification, (2) the use of parametrical and topological operators, and (3) an architecture-based competition scheme that sustains solution diversity and avoids premature convergence.

1. **Evolutionary complexification:** NEMODPS begins with a population of uniform simple networks, i.e., neural networks composed of just input and output layers, fully connected, with randomly initialized connection weights. As the evolution proceeds, neural architectures gradually complexify by including more architectural elements (nodes and connections) in the network’s hidden layer, which connects inputs to outputs. These elements are randomly generated by topological evolutionary operators and selected by evolutionary pressure.
2. **Parametrical and topological operators:** EAs use evolutionary operators such as mutation and crossover to recombine existing individual parameters to generate new individuals. NE evolves individual architectures along with their parameters, and therefore it includes both parametrical, and topological mutation and crossover. In particular, the topological mutation operator performs a randomized addition of a node (sigmoidal or Gaussian) or a connection to an individual. Topological crossover assigns the offspring a mix of the parents’ architectures. NEMODPS coordinates the topological and parametrical search in a dual timescale: parametrical mutation and crossover takes place every generation, while topological variations happen on a slower timescale, every few generations, to allow the competition scheme to protect solution diversity.
3. **Competition scheme:** at every generation, the population is divided into species of individuals with similar topologies. Species compete among each other for their ability to reproduce, so that a larger offspring is assigned to well performing ones. A fitness sharing mechanism penalizes numerous species preventing them from taking over the entire population causing loss of topological diversity and premature convergence. NEMODPS generalizes the fitness sharing strategy for MO problems, rewarding species with Pareto efficient individuals, and penalizing species whose individuals are located in crowded region of the objectives space in order to encourage the exploration of the entire tradeoff space.

2.3 Extraction of optimal decision from a Perfect Operating Policy

Following *Giuliani et al. (2015)*, the Perfect Operating Policy π^{POP} is designed by solving Problem 1 under the hypothesis of deterministic knowledge of the trajectory ε_1^H of external drivers over the entire evaluation horizon H at any given time step, $\pi^{POP} = \pi^{POP}(s_t, t, \varepsilon_1^H)$ and can be solved via various open loop deterministic control methods (examples can be found in, e.g., *Dobson et al., 2019; Macian-Sorribes and Pulido-Velazquez, 2020*). Here, we solve the problem with Deterministic Dynamic Programming (DDP). Such a deterministic policy can be considered the optimal reference for improving a basic policy design, but cannot be realistically implemented in a real-world system (e.g., *Denaro et al., 2017*). In order to obtain the trajectory of decision residuals e_t , we compare the decisions extracted from the π^{POP} with those extracted from the efficient policy π^R at a given round R , referring to the same state trajectory produced by the simulation of π^R . The difference in decisions extracted by the policy under design π^R , and the perfectly informed policy π^{POP} , is assumed to be due to their information gap. In a MO problem, π^R and π^{POP} are constituted by a set of Pareto efficient policies, therefore, each π^R policy is associated with the POP solution that displays the most similar tradeoff.

2.4 Termination criterion

SINEPS terminates at round $R > 1$ when the efficient Pareto set at Round R does not significantly dominate the Pareto set in the previous round: $\pi^R \not\prec \pi^{R-1}$, according to an appropriate metric. Several metrics could in principle be used to express dominance in a Pareto sense. Here, as suggested in *Giuliani et al. (2015)*, we use the hypervolume indicator (HV), which captures both the convergence of the Pareto front under examination \mathcal{F} to the optimal one \mathcal{F}^* , as well as the representation of the full extent of tradeoffs in the objective space. The hypervolume metric allows set-to-set evaluations, measuring the volume of objective space Y dominated (\preceq) by the considered approximate set. HV assumes values between 0 to 1, where Pareto fronts with higher HV are considered better. For this study, we consider the search terminated when the HV increase from round $R-1$ to round R is lower than 5%. Policies in round R are characterized by an increased complexity in the input layer, that however doesn't yield a significant performance increase. Therefore, round R is discarded, and the policies produced at round $R-1$ are considered final.

3 Case Study and Data

We consider the case study of Lake Como, a multipurpose regulated lake located in the southern Alpine belt, Italy (Fig. 2). The main tributary, and only emissary of the lake is the Adda river, whose sublacual reach originates in the southeastern branch of Lake Como, crosses the Po valley, and eventually serves as a tributary to the Po river downstream. In its course, part of its waters are withdrawn to irrigate four agricultural districts. The southwestern branch of Lake Como constitutes a dead end, and exposes the city of Como to flooding events. The Lake Como basin hydrological regime is snow-rainfall dominated, characterized by scarce winter and summer inflows, a large snowmelt peak in late spring, and a secondary rainfall peak in autumn.

The lake regulation has two conflicting aims of supplying water to downstream users by storing spring snowmelt peak, and minimizing flood risk on the lake shores by maintaining the lake level as low as possible, therefore, \mathbf{J} in eq. 1 is a bidimensional vector. On the basis of previous works (e.g., *Castelletti et al., 2010*), these two objectives are defined as:

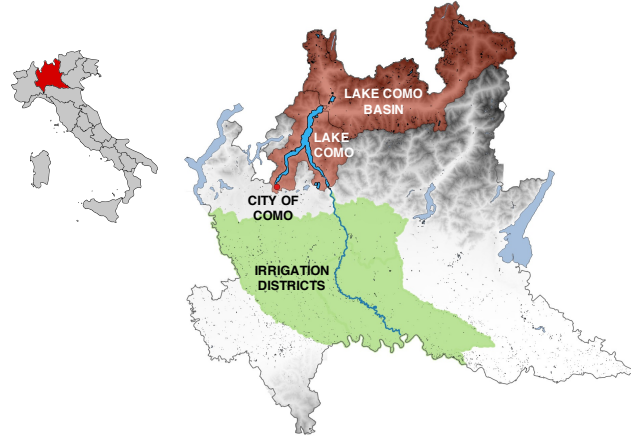


Figure 2. On the left, Lombardy region is highlighted in a map of Italy. On the right, a physical map of Lombardy, comprising Lake Como basin, in red, Lake Como, the city of Como, and the irrigation district downstream the lake.

Flood days: the average number of annual flood days, defined as days in which the lake level h_t is above the flood threshold $\bar{h} = 1.24$ m, i.e.,

$$J^{flood} = \frac{1}{N_y} \sum_{t=0}^{H-1} g_{t+1}^{flood}, \quad g_{t+1}^{flood} = \begin{cases} 1 & \text{if } h_{t+1} \geq \bar{h} \\ 0 & \text{if } h_{t+1} < \bar{h} \end{cases} \quad (3)$$

where N_y is the number of years in the simulation horizon.

Water supply deficit: the daily average squared water deficit with respect to the daily downstream demand w_t , subject to the minimum flow constraint $q^{MEF} = 5 \text{ m}^3/\text{s}$ to guarantee environmental stakes. Downstream demand is mainly driven by irrigation and is highest during the crop growing season of spring and summer. The quadratic formulation is selected with the aim of penalizing severe deficits in a single time step, while allowing for more frequent, small shortages. i.e.,

$$J^{irr} = \frac{1}{H} \sum_{t=0}^{H-1} (\max(w_t - (r_{t+1} - q^{MEF}), 0))^2 \quad (4)$$

The release decision is conditioned on an annual cyclostationary time index, and thus the decision at the end of the time horizon is no different than during the equivalent period of all previous years. For this application, we used Lake Como inflow data for a 10 year optimization horizon from 1997 to 2006 included. This time span contains a diverse range of hydrological conditions, including average and extreme years, from the 2005 record drought to the late 2000 high inflow pulses. Optimal policies are then tested on three validation chunks: an extended 20-years validation from 1977-1996, a combination of extreme dry years (1949, 1962, 1990, 1994, 2007), and wet years (1951, 1960, 1977, 2008, 2014) selected by searching the driest and wettest years from the available historical record of inflows to Lake Como (1947-2014), discarding the calibration years.

The set of candidate policy inputs employed in this analysis includes perfect forecasts of the lake inflow computed over the historical timeseries at different lead times, ranging from one day to over 6 months (Table 1). The forecasts are of two types: i) Cumulated inflows, which represent the cumulative inflows over a given lead time, and ii) Inflow Anomaly, which corresponds to the anomalies in inflow with respect to the inflow cyclostationary mean, cumulated over a given lead time. As argued in the introduction,

the aim of this methodological contribution is to demonstrate that the optimal policy representation varies with the objective tradeoff, and, therefore, one single policy representation is inadequate to represent the entire tradeoff space. The risk of using real forecasts in order to prove this concept is that the forecast bias may introduce noise and errors, and ultimately alter the information selection. Therefore, as per previous works (Zhao *et al.*, 2011; Denaro *et al.*, 2017), we made the modeling choice of using perfect forecasts with the aim of searching the optimal policy representation for the system, given its hydrology, physical characteristics, and objectives, and without being biased by errors in forecast products.

Lead time	Feature name	
	Cumulated inflow	Inflow Anomaly
1	C1	A1
2	C2	A2
3	C3	A3
5	C4	A5
7	C7	A7
14	C14	A14
21	C21	A21
28	C28	A28
51	C51	A51
62	C62	A62
75	C75	A75
90	C90	A90
120	C120	A120
145	C145	A145
200	C200	A200

Table 1. Dataset of candidate policy inputs comprising perfect inflow forecasts in terms of cumulated inflows and anomalies at various lead times.

3.1 Experimental Settings

SINEPS was run for 20 independently initialized and randomized seeds. In each seed, the termination criterion (described in Section 2.4) is met at the 4th round, which is responsible for no tangible advancement in the Pareto front, (lower than 5%), therefore, we retain as efficient solutions those generated at round 3. At each round, NEMODPS is run for a Number of Function Evaluations (NFE) equal to 600 thousands, with populations of 600 individuals. When new policy inputs are selected in step 4 of the methods, these are connected to the previously optimized policy architectures with an input-output connection. This set of individuals constitutes the initial population of the new round of NEMODPS optimization, in step 1 of round R2.

4 Results

4.1 Feature selection and policy design

Figure 3 reports the Pareto fronts resulting from 3 optimization rounds of SINEPS with respect to the two objectives of Water supply deficit (vertical axis) and Flood days (horizontal axis), both to be minimized as indicated by the arrows. The black square in the bottom left corner of the graph represents the ideal performance of the POP. In accordance to other studies on the same water system, we find that the conflicts between water supply and flood objectives in Lake Como disappear under the assumption of per-

507 perfect knowledge of future inflow (*Denaro et al.*, 2017). An operating policy with full fore-
508 sight is able to guarantee a sufficient flood pool to buffer the peak inflow and avoid over-
509 flow when physically possible, while storing in the lake any excess of water to be used
510 for irrigation purposes during the dry season. Therefore, the deterministic solution of
511 this MO problem does not yield a Pareto front of efficient solutions, but collapses to a
512 single optimal point into the objective space. However, in the absence of a perfect fu-
513 ture foresight, we expect that the addition of tailored information can reduce conflicts
514 between water users.

515 The first round of NEMODPS optimization, conditioned upon basic information
516 only, produces the Pareto front of white circles that lays in the top right portion of the
517 objective space in Figure 3a, showing a sharp conflict between the two operating objec-
518 tives. Additionally, a concavity can be recognized in the central region of the Pareto front,
519 for values of the Flood objective between 20 to 80. Concave regions of the front are usu-
520 ally regarded as disadvantageous tradeoffs, as one objectives degrades more than pro-
521 portionally to the second objective's improvement. The normalized HV indicator (panel
522 b) relative to round R1 scores 0.142, indicating a large space for improvement between
523 POP and R1.

524 Prior to the second NEMODPS optimization round, a feature selection routine iden-
525 tifies the most suitable variables to inform the operating policies via a two-step selec-
526 tion process. First, promising features are identified based on their correlation, measured
527 in SU, with the policy error trajectory, representative of its information gap (Figure 1,
528 box 3). Figure S2 of the Supplementary Information shows examples of error trajecto-
529 ries against the forecast anomaly lead time that scores the highest SU for different ob-
530 jective tradeoffs. Second, a population comprising all the promising features is evolved
531 via NEMODPS, and the fittest representations prevail through evolutionary competi-
532 tion (Figure 1, box 1 for $R > 1$). In particular, only a subset of the policy representations
533 preselected via SU is likely to survive the evolutionary selection pressure, meanwhile new
534 individuals are generated by recombining existing ones and enabling well-performing rep-
535 resentations to survive in future generations and establish in the final Pareto front. Fig-
536 ure S1 of the Supplementary Information reports the intermediate results of the two-fold
537 Feature Selection process, highlighting that evolutionary competition is key to identify
538 a contained and relevant feature set for policy representation.

539 The result of the second NEMODPS optimization round are represented in Fig-
540 ure 3a with colored triangles. The more informed policies significantly outperform R1,
541 scoring an over 3-fold increase in the HV metric. The color of the triangle corresponds
542 to the new feature added to the policy representation, and divides the R2 front in two,
543 around its middle and in correspondence to the persisting concavity in the Pareto front.
544 The analysis of the selected information may uncover unexpected results: flood-inclined
545 policies do not select short term predictions of fast inflow peaks, but long forecasts lead
546 times (75 days). Vice versa, water supply-inclined policies select, in comparison, slightly
547 shorter lead times (62 days) instead of preferring season-long look-ahead. This behav-
548 ior can be explained from the point of view of conflict mitigation. A minimally-represented
549 flood-inclined policy has, in fact, already developed a solid strategy to prevent floods when
550 physically possible, namely, keeping a low lake level for the most part of the year to al-
551 ways count on a buffer pool to accommodate incoming inflow peaks. This strategy is valid
552 from a lakeshore protection perspective, yet, comes at a remarkable price in terms of wa-
553 ter supply. Such policy, therefore, does not require any additional information on up-
554 coming inflow peaks, as the lake is virtually always ready to buffer them. On the con-
555 trary, it can significantly benefit from a longer term information on how to improve ir-
556 rigation while still remaining strongly flood risk-adverse, thereby alleviating water sup-
557 ply deficit downstream, and mitigating conflicts between water users. In fact, by com-
558 paring flood conservative policies of R1 and R2 (left region of the Pareto fronts), we no-
559 tice that the added information has the effect of improving the policies in the direction

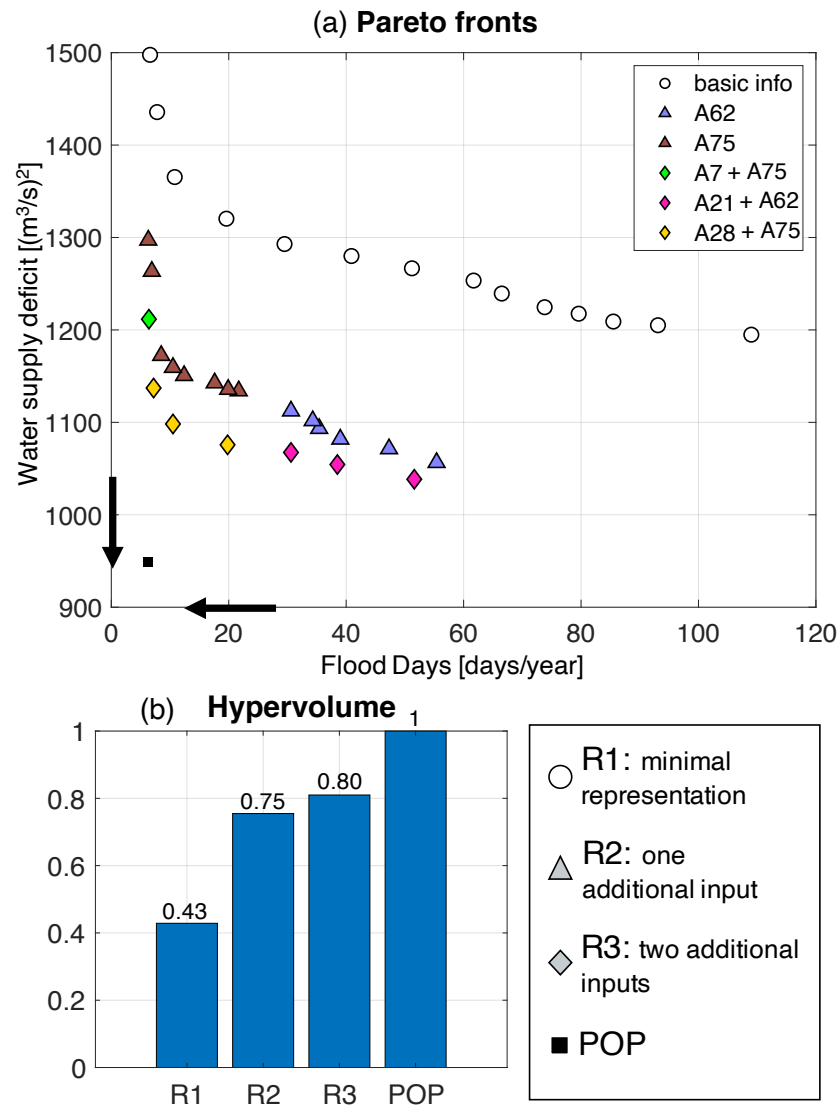


Figure 3. Panel (a): Performance obtained by different Lake Como operating policies with respect to the two cost objectives of water supply deficit (vertical axis) and Flood days (horizontal axis). The black square indicates the ideal performance of the POP, white circles the performance of efficient policies designed at round R1, triangles refer to policies at round R2, and diamonds at round R3. For rounds R2 and R3, the shape color is associated with the information added to the feature representation. Panel (b) shows the improvements in the Hypervolume indicator across different rounds, normalized to the value of hypervolume scored by POP.

560 of a significantly lower irrigation deficit, at no cost for the flood objective. The long lead
 561 time information selected by flood oriented policies is thus employed to minimize objec-
 562 tives conflicts, rather than further improve the flood objective. The other half of the Pareto
 563 front selects a shorter lead time, which allows policies to move both in the direction of
 564 a reduced flood and irrigation damage. Overall, however, this first round of information
 565 selection produces the largest improvement in the reduction of the water supply deficit
 566 by employing forecast with a long lead time (2 months or more). This selection is co-
 567 herent with the multi-seasonal nature of the water supply operations in a snow-dominated
 568 system like the one considered in this study. In particular, the reservoir is used to cre-
 569 ate the seasonal storage by impounding the spring snowmelt-driven inflow peak and dis-
 570 tribute it throughout the irrigation season, from spring to autumn, when water supply
 571 demand is highest. Forecast lead times of 2+ months are thus used to plan summer ir-
 572 rigation and inform the implementation of effective hedging rules when natural water
 573 availability does not meet demand. Lastly, policies select the anomaly in cumulated flow
 574 (A75, A62), rather than the flow cumulation, as it is a better indication of whether the
 575 system is entering a dry season and hedging strategies should be activated.

576 The third optimization round includes a second additional information in the pol-
 577 icy input set generating further improvement in the HV indicator. The Pareto front of
 578 round R3 not only dominates the fronts of the previous rounds, but also resolves their
 579 concavity generating a fully convex front, where it is possible to identify a knee. Con-
 580 trary to the previous round, the front shift between R2 and R3 is mainly horizontal, i.e.,
 581 contributing to a Flood objective improvement rather than an water supply improve-
 582 ment. Accordingly, the policy inputs selected in this round have a much shorter lead time,
 583 between 1 and 4 weeks. The solutions that at this round select the longer lead time, 4
 584 weeks, are those showing a diagonal improvement that unfolds in both objective direc-
 585 tions. We note that the by using perfect forecasts to inform the policies, the results shown
 586 in our work are upper bounds of what could be achievable with real forecasts in the sys-
 587 tem. For a demonstrative comparison of the performance using real forecasts instead of
 588 perfect forecasts, refer to section S4 of the SI.

589 It is worth noting that the optimal representations always select the anomaly in
 590 flow cumulation, over the flow cumulation. Cumulation time-series are analogous to their
 591 anomalies except for an additive cyclostationary, term which corresponds to the annual
 592 climatology and expresses the standard hydrological seasonality. However, the policy min-
 593 imal representation $\pi^{R1} = \pi^{R1}(d_t, s_t)$ already contains a cyclostationary time index
 594 d_t , which encapsules the climatology. As a consequence, it seems rational for the pol-
 595 icy to prefer the selection of an anomaly information over a partially redundant cumu-
 596 lative information. Additionally, it is common for medium-to-long term forecasts prod-
 597 ucts to produce forecast anomalies rather than cumulation (*Crochemore et al.*, 2020).

598 4.2 The role of information for conflict mitigation

599 In Figure 4 we explore how added information is employed by progressively informed
 600 policies for a given tradeoff. This analysis focuses on the solutions located along the lilac
 601 vertical line in panel (a), corresponding to an average of 6.3 flood days a year. This trade-
 602 off was chosen in order to compare the 4 Pareto fronts only in terms of the water sup-
 603 ply objective, for a given flood performance. A common cyclostationary behavior emerges
 604 for different policy representations in panel (b). The lake recharges in May, in correspon-
 605 dence to the onset of the irrigation season, reaches a level peak around late June, fol-
 606 lowed by an emptying phase lasting for the entire irrigation season until September/October,
 607 when abundant rains cause a new level increase. In the POP, perfect future foresight in-
 608 forms the policy on the exact onset of inflow peaks, allowing to timely generate an ad-
 609 equate flood pool to contain them, while keeping, on average, a high lake level that en-
 610 sures water availability to supply downstream irrigation demand. Whenever the full tra-
 611 jectory of future disturbance is not available, policies have to be more conservative to-

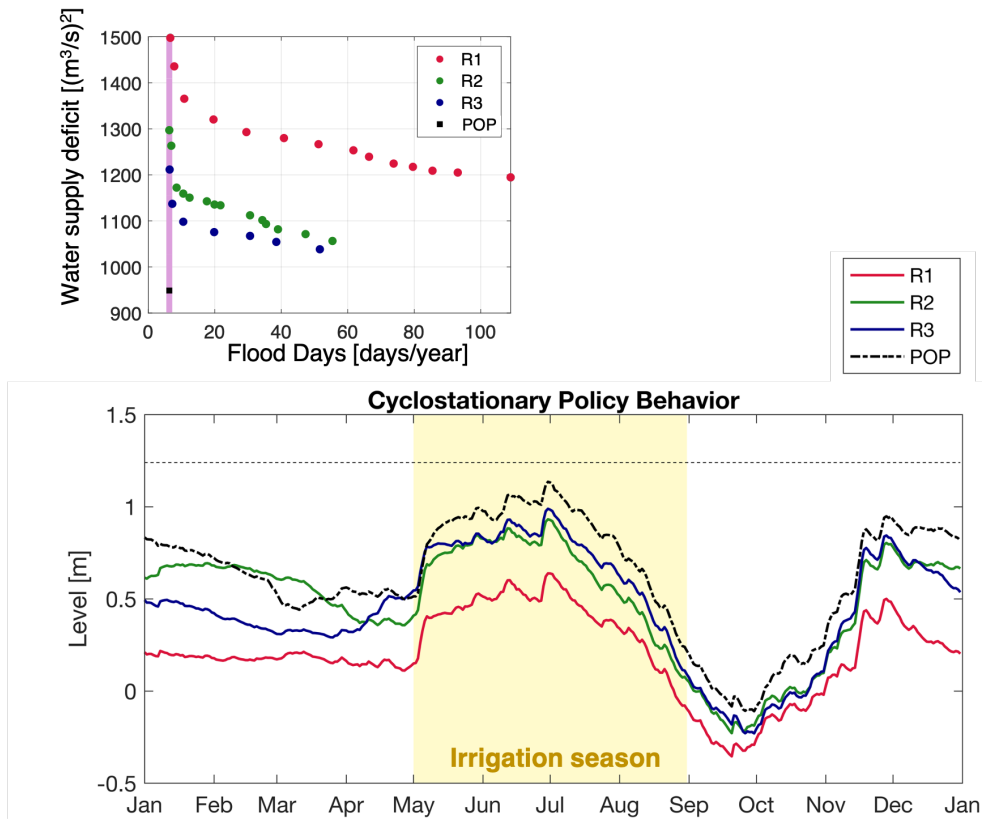


Figure 4. Cyclostationary behavior of efficient policies across different optimization rounds. The investigated policies are aligned along the lilac line in the Pareto front of panel (a) and yield an average number of flood days equal to 6.3, and different values with respect to the water supply objective. In panel (b), their cyclostationary behavior is shown, and contrasted with the Perfect Operating Policy.

612 wards flood events, thereby keeping a lower lake level to buffer possible incoming inflow
 613 peaks, at the expense irrigation availability. This behavior is sharper in the minimally
 614 informed round R1 (red line), while more informed policies can confidently maintain a
 615 fuller lake during the summer, resulting in a smaller water deficit downstream, without
 616 damaging the flood objective. Cyclostationary behaviors outside the irrigation season
 617 are fairly divergent, however, the system's winter downstream demand is almost negli-
 618 gible with respect to summer demand, thereby not contributing significantly to the wa-
 619 ter supply objective performance.

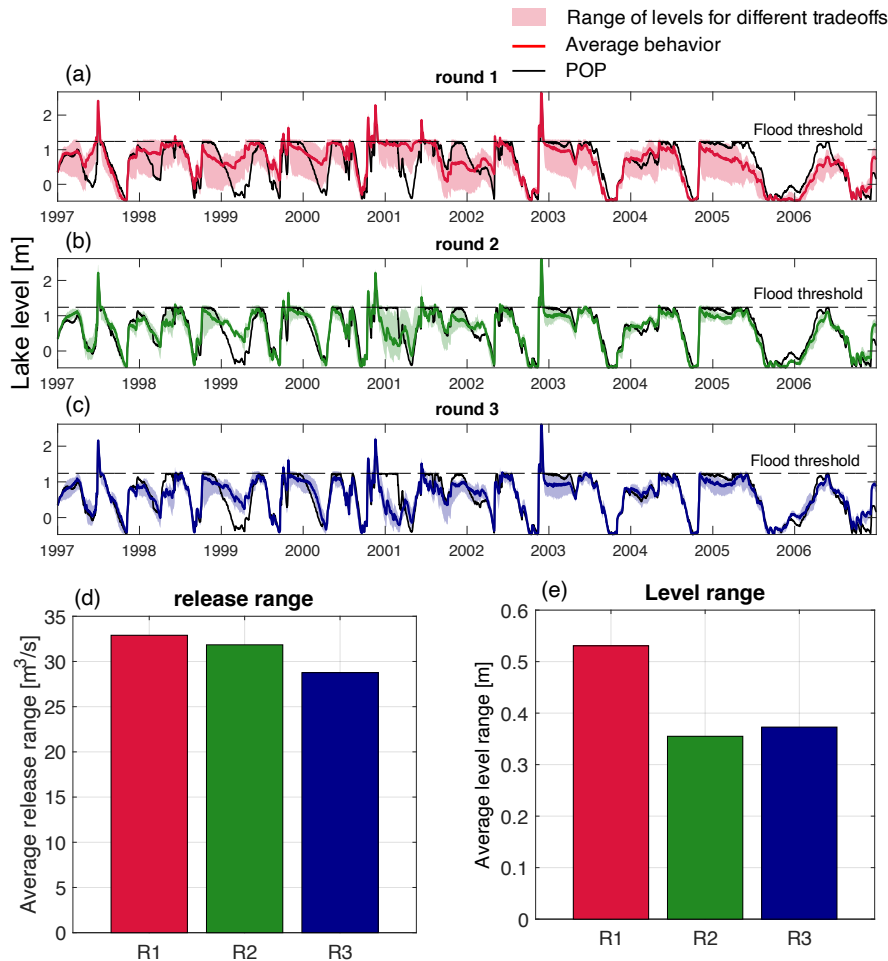


Figure 5. Conflict mitigation. Panels (a), (b), and (c) report the range of lake levels yielded by all the Pareto efficient policies designed at the given optimization round across different tradeoffs. The optimal trajectory is reported in every panel in black for reference. The average round-specific release range is quantified in the barplot of panel (d), while the lake level range is shown in panel (e).

620 In Figure 5, we analyze how a refinement in policy representation operationally mod-
 621 ifies lake regulation towards conflict mitigation. The shaded area in panels (a), (b), and
 622 (c) delimits the ensemble of lake level trajectories associated to the set of Pareto efficient
 623 policies produced in a given round, while the central colored bold line represents the av-
 624 erage behavior. The optimal POP trajectory is reported in black for reference. The wide-
 625 neness of the shaded area indicates the range of variability in operations spanned by the

626 efficient policies, where a thick area indicates that different tradeoffs are associated with
 627 diverse operations, and a narrow area suggests similar operations even across opposite
 628 tradeoffs. The plots show a visible narrowing in the operational variability from the first
 629 round to the following ones. Operationally, this translates into a mitigated conflict be-
 630 tween water users, as different interests tend to converge towards a common efficient pol-
 631 icy. This convergence is quantified in the barplots showing the average daily range in lev-
 632 els (panel e) and releases (panel d) associated to different policies in the Pareto set re-
 633 sulting from a given round. The addition of information in the policy representation shows
 634 a consistent reduction in release variability. Level variability significantly drops from round
 635 R1, where Lake Como is operated at an average difference of more than 53 cm for dif-
 636 ferent tradeoffs, to about 35 cm in round R2. R3 shows a slight increase in variability
 637 that is however below 2 cm, and can be considered negligible.

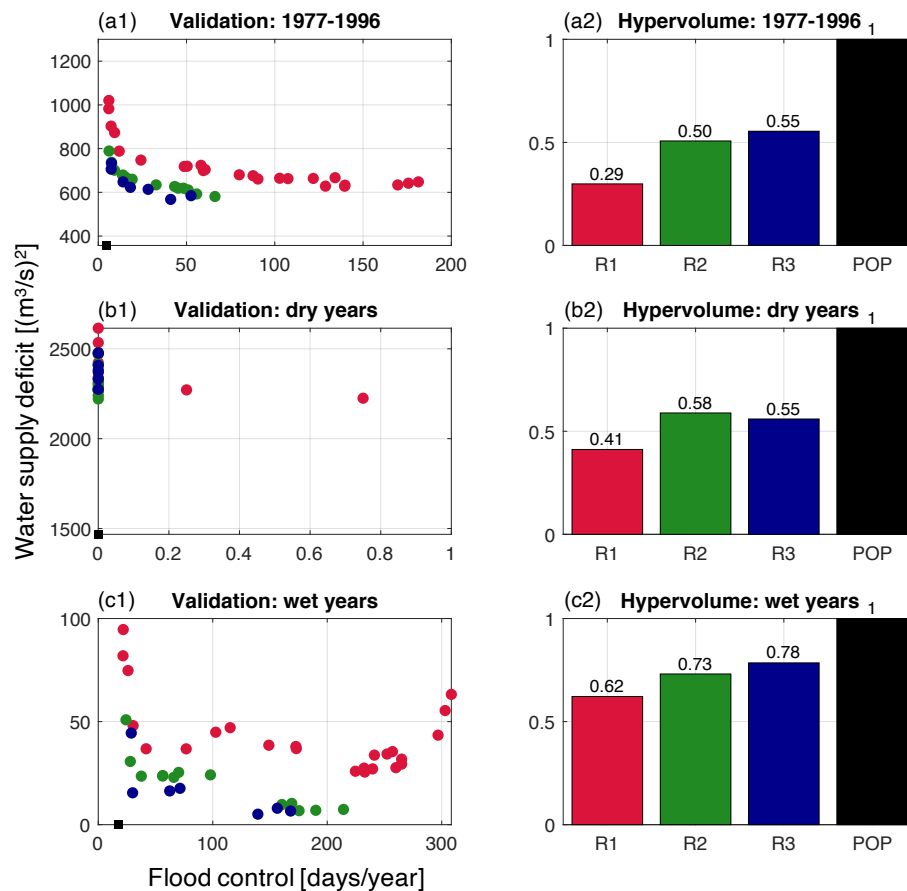


Figure 6. Validation of optimal policies for the three rounds of SINEPS for a 20-year evaluation horizon (panel a1 and a2), and two 5-year evaluation horizons composed of extreme dry (panels b1 and b2) and wet years (panels c1 and c2).

638 4.3 Policy validation in uncertain hydrological conditions

639 Figure 6 shows the re-evaluation of the optimal policies on three inflow trajectories,
 640 an extended 20-years horizon 1977-1996 (panel a1), an extreme dry (panel b1), and
 641 an extreme wet horizon (panels c1). Panels a2, b2, and c2 report the value of the HV in-

642 indicator computed for the different Rounds for the corresponding validation period. The
643 POP performance is reported for reference is each panel colored in black. The most in-
644 formed round R3 outperforms the other two in the 1977-1996 and wet-years datasets,
645 as quantified by the HV indicator and evident by the Pareto front of optimal validation
646 policies, which is composed by R3 solutions except for sporadic instances of R2 solutions
647 in panel c1. In the dry years dataset, one R2 solution achieve slightly lower water sup-
648 ply deficit compared to R3, but with a fairly negligible difference, under 3%. This anal-
649 ysis shows that the performance improvement resulting from an enhanced information
650 set persists in validation proving the robustness of the information selection technique
651 across highly diverse hydrological conditions.

652 5 Conclusions

653 In the past, reservoir operating rules were conditioned upon basic information sys-
654 tems comprising time index and reservoir storage (*Hejazi et al.*, 2008). However, the po-
655 tential of enhancing the performance of water system operations using information on
656 current or future water availability has long been recognized by researchers and prac-
657 titioners alike. Despite many features can contribute to operations to some extent, it is
658 in general unclear what is the most effective information set to condition a given water
659 system, for a given tradeoff.

660 Moreover, previous studies have generally overlooked how defining one single pol-
661 icy representation to characterize the entire tradeoff space of multi-purpose systems can
662 be insufficient. The coexistence of fast and slow process dynamics, and different vulner-
663 abilities requires the search of a tradeoff-tailored policy representation. In this work, we
664 demonstrate for the first time that one input set is inadequate to inform the entire Pareto
665 front of efficient policies that constitutes the solution to a multi-objective problem. In
666 fact, when the policy search routine is allowed to evolve heterogeneous input sets, the
667 selected optimal policy representation will vary Pareto-dynamically with the tradeoff.
668 In this work, we propose SINEPS, a novel framework for automatic, tradeoff-dynamic
669 feature representation and policy learning. SINEPS starts with a population of minimal
670 policies and gradually complexifies their feature representation by selecting variables that
671 surrogate the policy information deficit, measured by comparison to a Perfect Operat-
672 ing Policy. Policies' architectures are adjusted accordingly, in order to accomodate new
673 inputs and support more complex behaviors. We apply SINEPS to the case study of Lake
674 Como, characterized by conflicting heterogeneous objectives, and we use a dataset of de-
675 terministic inflow forecasts at different lead times as candidate policy inputs.

676 Results show that different objective tradeoffs benefit from different information
677 sets with unexpected, but insightful, outcomes. Flood-conservative policies select fore-
678 casts with long lead times, thereby improving water supply performance without increas-
679 ing flood failures. water supply-inclined policies select, in comparison, shorter lead times
680 achieving better flood and water supply results. Not only we notice a trend in the in-
681 formation selected for different tradeoffs, but also across subsequent selection rounds.
682 The first forecast included in the representation at the second round counts on a over
683 2 months-ahead lead time, and produces the largest improvement in the direction of a
684 lower water supply deficit, and only partially, flood mitigation. In round three, lead times
685 are shorter than a month, enhancing primarily flood mitigation skills. Overall, the search
686 for a tradeoff-specific feature representation demonstrates the potential to significantly
687 enhance the water system overall reliability, resilience towards both dry and wet extremes,
688 while reducing conflicts across conflicting water uses.

689 Lastly, it is important to note that policy representation in water resources man-
690 agement should not be considered a static concept, but should dynamically adapt in re-
691 sponse to variations in the ever-evolving boundary conditions that coupled human-natural
692 systems are exposed to. In particular, the optimal policy representation could change

693 in response to variations in socio-economic drivers e.g., a water user experiencing unprece-
 694 dented and more frequent failures; climatic drivers, i.e., an increased likelihood of one
 695 of more class of extreme event; and physical drivers, e.g., when a new water user or in-
 696 frastructure is included in the system. When one or more of these drivers change, the
 697 previous policy representation may not be adequate to represent the new system condi-
 698 tions and should be updated accordingly. The SINEPS framework can be run frequently
 699 to monitor and adapt to such changes with a rolling calibration horizon that includes
 700 new observations as they become available. A critical challenge yet to address is to de-
 701 termine when and how to timely update the feature representation by means of appro-
 702 priate triggers.

703 Data availability statement

704 The data used in this work are freely available upon request from Consorzio del-
 705 lAdda at <https://addaconsorzio.it/>.

706 References

- 707 Akrou, R., M. Schoenauer, and M. Sebag (2012), April: Active preference learning-
 708 based reinforcement learning, in *Joint European Conference on Machine Learning
 709 and Knowledge Discovery in Databases*, pp. 116–131, Springer.
- 710 Alvernaz, S., and J. Togelius (2017), Autoencoder-augmented neuroevolution for
 711 visual doom playing, in *2017 IEEE Conference on Computational Intelligence and
 712 Games (CIG)*, pp. 1–8, IEEE.
- 713 Anghileri, D., N. Voisin, A. Castelletti, F. Pianosi, B. Nijssen, and D. Lettenmaier
 714 (2016), *Value of long-term streamflow forecasts to reservoir operations for water
 715 supply in snow-dominated river catchments*, *Water Resources Research*.
- 716 Assael, J.-A. M., N. Wahlström, T. B. Schön, and M. P. Deisenroth (2015), Data-
 717 efficient learning of feedback policies from image pixels using deep dynamical
 718 models, *arXiv preprint arXiv:1510.02173*.
- 719 Bengio, Y., A. Courville, and P. Vincent (2013), Representation learning: A review
 720 and new perspectives, *IEEE transactions on pattern analysis and machine intelli-
 721 gence*, 35(8), 1798–1828.
- 722 Block, P. (2011), Tailoring seasonal climate forecasts for hydropower operations.,
 723 *Hydrology & Earth System Sciences*, 15(4).
- 724 Blum, A. L., and P. Langley (1997), Selection of relevant features and examples in
 725 machine learning, *Artificial intelligence*, 97(1-2), 245–271.
- 726 Castelletti, A., S. Galelli, M. Restelli, and R. Soncini-Sessa (2010), Tree-based re-
 727 inforcement learning for optimal water reservoir operation, *Water Resources
 728 Research*, 46(9).
- 729 Castelletti, A., R. Fedorov, P. Fraternali, and M. Giuliani (2016), Multimedia on
 730 the mountaintop: Using public snow images to improve water systems opera-
 731 tion, in *Proceedings of the 24th ACM international conference on Multimedia*, pp.
 732 948–957.
- 733 Crochemore, L., M.-H. Ramos, and I. Pechlivanidis (2020), Can continental mod-
 734 els convey useful seasonal hydrologic information at the catchment scale?, *Water
 735 Resources Research*, 56(2), e2019WR025,700.
- 736 Cunningham, P. (2008), Dimension reduction, in *Machine learning techniques for
 737 multimedia*, pp. 91–112, Springer.
- 738 Curran, W., T. Brys, D. Aha, M. Taylor, and W. D. Smart (2016), Dimensionality
 739 reduced reinforcement learning for assistive robots, in *2016 AAAI Fall Symposium
 740 Series*.
- 741 Denaro, S., D. Anghileri, M. Giuliani, and A. Castelletti (2017), Informing the
 742 operations of water reservoirs over multiple temporal scales by direct use of hydro-

- 743 meteorological data, *Advances in water resources*, 103, 51–63.
- 744 Desreumaux, Q., P. Côté, and R. Leconte (2014), Role of hydrologic information in
745 stochastic dynamic programming: a case study of the kemano hydropower system
746 in british columbia, *Canadian Journal of Civil Engineering*, 41(9), 839–844.
- 747 Dobson, B., T. Wagener, and F. Pianosi (2019), An argument-driven classification
748 and comparison of reservoir operation optimization methods, *Advances in Water*
749 *Resources*, 128, 74–86.
- 750 Doering, K., J. Quinn, P. M. Reed, and S. Steinschneider (2021), Diagnosing the
751 time-varying value of forecasts in multiobjective reservoir control, *Journal of Wa-*
752 *ter Resources Planning and Management*, 147(7), 04021,031.
- 753 Fletcher, S., M. Lickley, and K. Strzepek (2019), Learning about climate change
754 uncertainty enables flexible water infrastructure planning, *Nature communications*,
755 10(1), 1–11.
- 756 Floreano, D., P. Dürr, and C. Mattiussi (2008), Neuroevolution: from architectures
757 to learning, *Evolutionary Intelligence*, 1(1), 47–62.
- 758 Gal, S. (1979), Optimal management of a multireservoir water supply system, *Water*
759 *Resources Research*, 15(4), 737–749.
- 760 Gaudel, R., and M. Sebag (2010), Feature Selection as a One-Player Game, in *In-*
761 *ternational Conference on Machine Learning*, ICML 2010 Conference Proceedings
762 Book, pp. 359–366, Haifa, Israel.
- 763 Giuliani, M., and A. Castelletti (2019), Data-driven control of water reservoirs using
764 el niño southern oscillation indexes, in *2019 IEEE International Conference on*
765 *Environment and Electrical Engineering and 2019 IEEE Industrial and Commer-*
766 *cial Power Systems Europe (EEEIC/I&CPS Europe)*, pp. 1–5, IEEE.
- 767 Giuliani, M., J. Herman, A. Castelletti, and P. Reed (2014), Many-objective reser-
768 voir policy identification and refinement to reduce policy inertia and myopia in
769 water management, *Water Resources Research*, 50(4), 3355–3377.
- 770 Giuliani, M., F. Pianosi, and A. Castelletti (2015), Making the most of data: an
771 information selection and assessment framework to improve water systems opera-
772 tions, *Water Resources Research*, 51(11), 9073–9093.
- 773 Giuliani, M., A. Castelletti, R. Fedorov, and P. Fraternali (2016a), Using crowd-
774 sourced web content for informing water systems operations in snow-dominated
775 catchments, *Hydrology and Earth System Sciences*, 10(5194), 20–5049.
- 776 Giuliani, M., A. Castelletti, F. Pianosi, E. Mason, and P. Reed (2016b), Curses,
777 tradeoffs, and scalable management: Advancing evolutionary multiobjective direct
778 policy search to improve water reservoir operations, *Journal of Water Resources*
779 *Planning and Management*, 142(2), 04015,050.
- 780 Giuliani, M., J. D. Quinn, J. D. Herman, A. Castelletti, and P. M. Reed (2018),
781 Scalable multiobjective control for large-scale water resources systems under un-
782 certainty, *IEEE Transactions on Control Systems Technology*, 26(4), 1492–1499.
- 783 Giuliani, M., M. Zaniolo, A. Castelletti, G. Davoli, and P. Block (2019), Detect-
784 ing the state of the climate system via artificial intelligence to improve seasonal
785 forecasts and inform reservoir operations, *Water Resources Research*, 55(11),
786 9133–9147.
- 787 Gleick, P. H. (2003), Global freshwater resources: soft-path solutions for the 21st
788 century, *Science*, 302(5650), 1524–1528.
- 789 Hachiya, H., and M. Sugiyama (2010), Feature selection for reinforcement learning:
790 Evaluating implicit state-reward dependency via conditional mutual information,
791 in *Joint European Conference on Machine Learning and Knowledge Discovery in*
792 *Databases*, pp. 474–489, Springer.
- 793 Hamlet, A. F., D. Huppert, and D. P. Lettenmaier (2002), Economic value of long-
794 lead streamflow forecasts for columbia river hydropower, *Journal of Water Re-*
795 *sources Planning and Management*, 128(2), 91–101.

- 796 Hejazi, M. I., and X. Cai (2011), Building more realistic reservoir optimization mod-
 797 els using data mining—a case study of shelbyville reservoir, *Advances in water*
 798 *resources*, *34*(6), 701–717.
- 799 Hejazi, M. I., X. Cai, and B. L. Ruddell (2008), The role of hydrologic informa-
 800 tion in reservoir operation—learning from historical releases, *Advances in water*
 801 *resources*, *31*(12), 1636–1650.
- 802 Herman, J. D., J. D. Quinn, S. Steinschneider, M. Giuliani, and S. Fletcher (2020),
 803 Climate adaptation as a control problem: Review and perspectives on dynamic
 804 water resources planning under uncertainty, *Water Resources Research*, *56*(2),
 805 e24,389.
- 806 Hundecha, Y., B. Arheimer, C. Donnelly, and I. Pechlivanidis (2016), A regional
 807 parameter estimation scheme for a pan-european multi-basin model, *Journal of*
 808 *Hydrology: Regional Studies*, *6*, 90–111.
- 809 Kroon, M., and S. Whiteson (2009), Automatic feature selection for model-based
 810 reinforcement learning in factored mdps, in *2009 International Conference on*
 811 *Machine Learning and Applications*, pp. 324–330, IEEE.
- 812 Lesort, T., N. Díaz-Rodríguez, J.-F. Goudou, and D. Filliat (2018), State representa-
 813 tion learning for control: An overview, *Neural Networks*, *108*, 379–392.
- 814 Libisch-Lehner, C., H. Nguyen, R. Taormina, H. Nachtnebel, and S. Galelli (2019),
 815 On the value of enso state for urban water supply system operators: Opportuni-
 816 ties, trade-offs, and challenges, *Water Resources Research*, *55*(4), 2856–2875.
- 817 Liu, D.-R., H.-L. Li, and D. Wang (2015), Feature selection and feature learning for
 818 high-dimensional batch reinforcement learning: A survey, *International Journal of*
 819 *Automation and Computing*, *12*(3), 229–242.
- 820 Loscalzo, S., R. Wright, and L. Yu (2015), Predictive feature selection for genetic
 821 policy search, *Autonomous Agents and Multi-Agent Systems*, *29*(5), 754–786.
- 822 Macian-Sorribes, H., and M. Pulido-Velazquez (2020), Inferring efficient operating
 823 rules in multireservoir water resource systems: A review, *Wiley Interdisciplinary*
 824 *Reviews: Water*, *7*(1), e1400.
- 825 MacKay, D. J. (2003), *Information theory, inference and learning algorithms*, Cam-
 826 bridge university press.
- 827 Maidment, D. R., and V. T. Chow (1981), Stochastic state variable dynamic pro-
 828 gramming for reservoir systems analysis, *Water Resources Research*, *17*(6), 1578–
 829 1584.
- 830 Molteni, F., T. Stockdale, M. Balmaseda, G. Balsamo, R. Buizza, L. Ferranti,
 831 L. Magnusson, K. Mogensen, T. Palmer, and F. Vitart (2011), *The new ECMWF*
 832 *seasonal forecast system (System 4)*, vol. 49, European Centre for Medium-Range
 833 Weather Forecasts Reading.
- 834 Morimoto, J., S.-H. Hyon, C. G. Atkeson, and G. Cheng (2008), Low-dimensional
 835 feature extraction for humanoid locomotion using kernel dimension reduction, in
 836 *2008 IEEE International Conference on Robotics and Automation*, pp. 2711–2716,
 837 IEEE.
- 838 Munk, J., J. Kober, and R. Babuška (2016), Learning state representation for deep
 839 actor-critic control, in *2016 IEEE 55th Conference on Decision and Control*
 840 *(CDC)*, pp. 4667–4673, IEEE.
- 841 Nouri, A., and M. L. Littman (2010), Dimension reduction and its application to
 842 model-based exploration in continuous spaces, *Machine Learning*, *81*(1), 85–98.
- 843 Oh, J., S. Singh, and H. Lee (2017), Value prediction network, in *Advances in Neural*
 844 *Information Processing Systems*, pp. 6118–6128.
- 845 Pechlivanidis, I., L. Crochemore, J. Rosberg, and T. Bosshard (2020), What are
 846 the key drivers controlling the quality of seasonal streamflow forecasts?, *Water*
 847 *Resources Research*, *56*(6), e2019WR026,987.
- 848 Quinn, J. D., P. M. Reed, M. Giuliani, and A. Castelletti (2017), Rival framings: A
 849 framework for discovering how problem formulation uncertainties shape risk man-

- agement trade-offs in water resources systems, *Water Resources Research*, 53(8), 7208–7233.
- Quinn, J. D., P. M. Reed, M. Giuliani, A. Castelletti, J. W. Oyster, and R. E. Nicholas (2018), Exploring how changing monsoonal dynamics and human pressures challenge multireservoir management for flood protection, hydropower production, and agricultural water supply, *Water Resources Research*, 54(7), 4638–4662.
- Quinn, J. D., P. M. Reed, M. Giuliani, and A. Castelletti (2019), What is controlling our control rules? opening the black box of multireservoir operating policies using time-varying sensitivity analysis, *Water Resources Research*, 55(7), 5962–5984.
- Shannon, C. E. (1948), A mathematical theory of communication, *The Bell system technical journal*, 27(3), 379–423.
- Si, W., J. Li, P. Ding, and R. Rao (2017), A multi-objective deep reinforcement learning approach for stock index future’s intraday trading, in *2017 10th International symposium on computational intelligence and design (ISCID)*, vol. 2, pp. 431–436, IEEE.
- Stanley, K. O., and R. Miikkulainen (2002), Efficient reinforcement learning through evolving neural network topologies, in *Proceedings of the 4th Annual Conference on Genetic and Evolutionary Computation*, pp. 569–577, Morgan Kaufmann Publishers Inc.
- Stanley, K. O., and R. Miikkulainen (2003), A taxonomy for artificial embryogeny, *Artif. Life*, 9(2), 93–130.
- Sturtevant, N. R., and A. M. White (2006), Feature construction for reinforcement learning in hearts, in *International Conference on Computers and Games*, pp. 122–134, Springer.
- Tan, M., R. Deklerck, J. Cornelis, and B. Jansen (2013), Phased searching with neat in a time-scaled framework: experiments on a computer-aided detection system for lung nodules, *Artificial intelligence in medicine*, 59(3), 157–167.
- Turner, S., and S. Galelli (2016), Regime-shifting streamflow processes: Implications for water supply reservoir operations, *Water Resources Research*, 52(5), 3984–4002.
- Turner, S., W. Xu, and N. Voisin (2019), Inferred inflow forecast horizons guiding reservoir release decisions across the united states, *Hydrology and Earth System Sciences Discussions*, pp. 1–25.
- Van Hoof, H., N. Chen, M. Karl, P. van der Smagt, and J. Peters (2016), Stable reinforcement learning with autoencoders for tactile and visual data, in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 3928–3934, IEEE.
- Weedon, G. P., G. Balsamo, N. Bellouin, S. Gomes, M. J. Best, and P. Viterbo (2014), The wfdei meteorological forcing data set: Watch forcing data methodology applied to era-interim reanalysis data, *Water Resources Research*, 50(9), 7505–7514.
- Whiteson, S., P. Stone, K. O. Stanley, R. Miikkulainen, and N. Kohl (2005), Automatic feature selection in neuroevolution, in *Proceedings of the 7th annual conference on Genetic and evolutionary computation*, pp. 1225–1232.
- Wright, R., S. Loscalzo, and L. Yu (2012), Embedded incremental feature selection for reinforcement learning, *Tech. rep.*, AIR FORCE RESEARCH LAB ROME NY INFORMATION DIRECTORATE.
- Xu, W., C. Zhang, Y. Peng, G. Fu, and H. Zhou (2014), A two stage bayesian stochastic optimization model for cascaded hydropower systems considering varying uncertainty of flow forecasts, *Water Resources Research*, 50(12), 9267–9286.
- Yang, W., J. Andréasson, L. Phil Graham, J. Olsson, J. Rosberg, and F. Wetterhall (2010), Distribution-based scaling to improve usability of regional climate model projections for hydrological climate change impacts studies, *Hydrology Research*,

904 41(3-4), 211–229.

905 Zaniolo, M., M. Giuliani, A. F. Castelletti, and M. Pulido-Velazquez (2018), Auto-
906 matic design of basin-specific drought indexes for highly regulated water systems,
907 *Hydrology and Earth System Sciences*, 22(4), 2409–2424.

908 Zaniolo, M., M. Giuliani, and A. Castelletti (2019), Data-driven modeling and con-
909 trol of droughts, *IFAC-PapersOnLine*, 52(23), 54–60.

910 Zaniolo, M., M. Giuliani, S. Sinclair, P. Burlando, and A. Castelletti (2021), When
911 timing matters: misdesigned dam filling impacts hydropower sustainability, *Nature*
912 *communications* .

913 Zaniolo, M., M. Giuliani, and A. Castelletti (2021), Neuro-evolutionary direct policy
914 search for multiobjective optimal control, *IEEE Transactions on Neural Networks*
915 *and Learning Systems*.

916 Zatarain, J. S., P. M. Reed, J. D. Quinn, M. Giuliani, and A. Castelletti (2017),
917 Balancing exploration, uncertainty and computational demands in many objective
918 reservoir optimization, *Advances in Water Resources*, 109, 196–210.

919 Zhang, L., and X. Chen (2021), Feature selection methods based on symmetric un-
920 certainty coefficients and independent classification information, *IEEE Access*, 9,
921 13,845–13,856.

922 Zhang, T. (2009), Adaptive forward-backward greedy algorithm for sparse learning
923 with linear models, in *Advances in Neural Information Processing Systems*, pp.
924 1921–1928.

925 Zhao, Q., X. Cai, and Y. Li (2019), Determining inflow forecast horizon for reservoir
926 operation, *Water Resources Research*, 55(5), 4066–4081.

927 Zhao, T., X. Cai, and D. Yang (2011), Effect of streamflow forecast uncertainty on
928 real-time reservoir operation, *Advances in water resources*, 34(4), 495–504.

929 Zhao, T., J. Zhao, J. R. Lund, and D. Yang (2014), Optimal hedging rules for reser-
930 voir flood operation from forecast uncertainties, *Journal of Water Resources*
931 *Planning and Management*, 140(12), 04014,041.