

WIDE-AREA DAMPING CONTROLLER VIA REINFORCEMENT LEARNING FOR POWER NETWORKS WITH WIND AND SOLAR FARMS

Muhammad Nadeem, MirSaleh Bahavarnia, and Ahmad F. Taha

Abstract—Accurately predicting power system behavior is becoming more complex with the increased penetration of uncertain wind and solar-based renewable resources. Hence, there is a growing motivation to transition from model-based feedback control strategies to completely model-free counterparts. Reinforcement learning (RL) is a key methodology in designing a model-free controller. Various studies have been carried out to study voltage/frequency control strategies via RL. However, they usually consider a simplified power system model either by completely neglecting system differential equations (and thus only modeling the system via power balance equations) or considering simplified generator models. Furthermore, damping system-wide oscillations after large disturbances are usually ignored in the controller design. In contrast, we propose an RL-based wide-area damping controller (WADC) for an advanced power system model with comprehensive higher-order generator dynamics, power electronics-based wind and solar models, and composite load dynamics. The presented controller sends control signals to synchronous generators, wind, and solar power plants to actively adjust their power and voltage setpoints—thereby providing damping to the system oscillations after large disturbances. Case studies demonstrate that the system’s transient stability can be significantly improved after large disturbances.

Index Terms—Reinforcement learning, solar and wind-based power system models, feedback control, adaptive control.

I. INTRODUCTION AND MOTIVATION

MODEL-based stability analysis and state/output feedback control have been the mainstream approaches to studying power system dynamics [1], [2]. These models, based on complex differential algebraic equations (DAE) offer invaluable insights into grid behavior, allowing for the design of sophisticated controllers that can effectively regulate voltage, frequency, and other critical parameters. However, the effectiveness of model-based control is contingent upon the accuracy of these models, which are inherently limited by simplifications and assumptions.

As power systems become more uncertain and dynamic due to the integration of renewables and the proliferation of distributed energy resources, accurately modeling their behavior becomes an increasingly daunting task. This leads to control strategies that are inadequate for real-world challenges [3]. In response to these evolving issues, there is a growing motivation to transition from traditional model-based feedback control to completely model-free control approaches [4]–[6]. This has also been highly encouraged by the recent

developments in the wide-area monitoring systems technologies, such as phasor measurement units (PMUs) and robust state estimation algorithms [4].

Using measurements received from the PMUs, these estimation techniques can efficiently estimate all the states of the power system including the states of generators, renewables, loads, and network. This paves the way for completely model-free control techniques in the future power grid.

RL is one of the key techniques in designing model-free feedback control strategies. This is because RL is a self-learning technique in which the agent (or in our case the controller itself) learns the optimal control policy dynamically by merely continuously interacting with its environment or simulation while satisfying the objective function or goal. In this regard, many researchers have recently proposed various RL control algorithms for power systems. In [7], [8], an RL-based optimal frequency controller has been proposed to minimize frequency deviations and the required control action for a simplified second-order power system model. Similarly, in [9] a voltage control strategy has been proposed using deep RL. The RL algorithm is trained to minimize the voltage deviation on all buses by actively adjusting generator bus voltages. However, system dynamic equations (i.e., generator model and power electronics-based model of renewables) are not considered; the corresponding power system is only modeled via power balance equations between generators and the network.

In [10]–[12] a multi-area automatic generation control (AGC) has been designed for frequency regulation using various RL-based techniques such as a deep deterministic policy gradient (DDPG), integral RL (actor-critic), and Q-learning based approaches. Similarly, in [13], [14] various reactive power control schemes have been proposed to improve system voltages via different RL-based techniques. A comprehensive review of the application of RL in the control and stability of power systems can be found in [3], [15].

All of the techniques in the literature usually consider simplified power system models (modeling lower order generator models), neglect power system algebraic constraints (current/power balance equations), and do not provide a control mechanism for (power-electronics)-based wind/solar dynamic models. Furthermore, most of the studies usually focus only on minimizing voltage/frequency deviations and do not consider damping system-wide low-frequency oscillations (LFOs) and ultra-low frequency oscillations (ULFOs) in their architecture. Damping LFOs and ULFOs are critical in the power grid as they can lead to system instability and limit power transfer capability. Furthermore, considering power system algebraic

The authors are with the Civil and Environmental Engineering Department, Vanderbilt University, 2201 West End Ave, Nashville, TN 37235, USA. Email addresses: muhammad.nadeem@vanderbilt.edu, mirsaleh.bahavarnia@vanderbilt.edu, ahmad.taha@vanderbilt.edu. This work is supported by the National Science Foundation under Grants 2152450 and 2151571.

constraints while designing a feedback control strategy is crucial because these constraints model the network dynamics. Thus without considering these constraints, the power system models cannot capture the effects of various load dynamics (such as constant impedance or motor-based loads) and topological changes (such as the tripping of transmission lines, etc) [2].

Some recent work has been carried out to design various centralized and decentralized RL-based damping controllers. In particular, the study [16] designs a wide-area damping controller based on a reduced system model. In [17] a dense centralized damping controller has been proposed based on a policy iteration algorithm. Similarly, in [18] a decentralized off-policy iteration-based damping controller has been designed. However, in all these studies again a simplified power system model is considered (considering lower order generator dynamics and neglecting system algebraic constraints) and no feedback control mechanism for renewables such as solar and wind-based power plants has been designed. To that end, we design an RL-based wide-area damping controller (RL-WADC) for an interconnected power system model with comprehensive 9th-order synchronous generator dynamics, detailed (power-electronics)-based wind and solar dynamics models, and composite load dynamics consisting of constant power, constant impedance, and motor loads. This brings the control problem closer to being more applicable and realistic. The technical contributions of the paper are as follows:

- We present a completely model-free approach to designing a wide-area damping controller for an advanced interconnected power system model. The presented approach considers: higher-order synchronous generators dynamics, detailed (power electronic)-based models of wind/solar-based power plants, models of the system algebraic constraints, and various load dynamics (constant power/impedance and motor-based loads) in its control architecture (Section II).
- The presented approach acts as a secondary control loop and sends additional control signals not only to the conventional power plants but also to the renewable resources (solar and wind-based power plants) so that they can participate in improving system transient stability after large disturbances in load/renewables (Sections III and IV).
- We test the performance of the presented RL-WADC under various load and renewable uncertainties on a modified IEEE 9-bus system—which is one of the widely used test systems in control/stability studies of power systems. The advantages of presented RL-WADC are also demonstrated by comparing the system dynamic response with conventional/primary controllers of the power system and with RL-WADC acting on top of them (Section V).

Paper Notation. All matrices and vectors are denoted using boldface type. Sets are represented using calligraphic font, such as \mathcal{R}, \mathcal{N} , etc. The symbol \mathbb{R}^a represents the set of column vectors with ‘a’ elements. Similarly, $\mathbb{R}^{a \times b}$ denotes the set

of real matrices with dimensions ‘a’ by ‘b’. The symbol \mathbf{O} represents a zero matrix, \mathbf{I} denotes an identity matrix of appropriate size, \cup denotes the union between two sets while the symbol \otimes represents the Kronecker product. In addition, $\mathbb{S}_{++}^{a \times b}$ signifies the set of positive definite matrices of size ‘a’ by ‘b’. The notation $\text{vec}(\mathbf{M})$ represents the vectorization of matrix \mathbf{M} which is computed by stacking the columns of \mathbf{M} on top of each other while $\bar{\mathbf{A}}$ denotes the half vectorization of a symmetric matrix \mathbf{A} with off-diagonal element taken as $2A_{ij}$. Similarly the notation $\bar{\mathbf{B}}$ represent the half vectorization plus the off-diagonal element of matrix \mathbf{B} , such as for $\mathbf{B} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$, $\bar{\mathbf{B}} = [1, 3, 4]^\top + [0, 2, 4]^\top$. Furthermore, the units for all the quantities are in per unit (p.u) unless otherwise specified, also, in some cases for brevity, the time stamp for some of the vectors is removed such as $\mathbf{u}_G(t)$ is written just as \mathbf{u}_G vice versa.

II. INTERCONNECTED POWER SYSTEM MODEL WITH WIND, SOLAR PLANTS, AND COMPOSITE LOADS

We model a power network with R solar farms, G conventional power plants, W wind-based power plants, and L_z, L_p, L_k , number of constant power, constant impedance, and motor-based loads, respectively. The complete electrical grid is represented via graph with \mathcal{E} set of transmission lines, and $\mathcal{N} = \{1, \dots, N\}$ set of nodes or buses. The set of buses \mathcal{N} contains; $\mathcal{G} = \{1, \dots, G\}$ set of buses with synchronous generators, $\mathcal{R} = \{1, \dots, R\}$ set of buses with solar farms, $\mathcal{W} = \{1, \dots, W\}$ set of buses connected to a wind-based power plant, \mathcal{L} set of buses collecting L_z, L_p , and L_k loads, and \mathcal{U} set of non-unit buses (meaning buses not connected to any elements), such that $\mathcal{N} = \mathcal{R} \cup \mathcal{W} \cup \mathcal{G} \cup \mathcal{U} \cup \mathcal{L}$.

The overall power network is described via nonlinear differential-algebraic equations (NDAEs) as follows:

$$\text{differential equations: } \dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}_a, \mathbf{x}_d, \mathbf{u}) \quad (1a)$$

$$\text{algebraic equations: } \mathbf{0} = \mathbf{h}(\mathbf{x}_a, \mathbf{x}_d, \mathbf{u}). \quad (1b)$$

The differential equations model the detailed dynamics of solar plants, wind-based generators, conventional power plants, and composite loads while the algebraic equations describe the algebraic constraints in the power network which are the power/current balance equations. In (1) $\mathbf{x}_a \in \mathbb{R}^{n_a}$ denotes algebraic states, $\mathbf{x}_d \in \mathbb{R}^{n_d}$ represents dynamic states and it lumps solar plants, wind-based plant, loads, and conventional power plant state variables, and $\mathbf{u} \in \mathbb{R}^{n_u}$ denotes the power systems control inputs. The more detailed explanations about each of these variables can be found in [1], [19], [20].

That being said, for control purposes we approximate the overall power system model in the following linear DAE format:

$$\mathbf{E}\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u} \quad (2)$$

where $\mathbf{x}(t) = [\mathbf{x}_d^\top \quad \mathbf{x}_a^\top]^\top \in \mathbb{R}^n$ is the overall state vector and \mathbf{E} is a singular matrix encoding algebraic constraints with rows of zeros.

III. PRELIMINARIES AND PROBLEM FORMULATION

Here we discuss the preliminaries and the problem formulation for designing a state feedback controller for the

interconnected power system model described in (2). The primary objective is to design a feedback control law that minimizes a certain performance criterion or cost function. To that end, let us consider the closed-loop dynamics of (2) along with its performance criterion vector $z \in \mathbb{R}^n$ as follows:

$$E\dot{x} = Ax + Bu_C \quad (3a)$$

$$z = Cx + Du_C \quad (3b)$$

where the control policy u_C for every k^{th} dispatch period is defined as

Control policy: $u_C(t) = u_0^k - K(x(t) - x^k)$

(4)

in which x^k is the state vector before fault/disturbance, u_0^k is the initial set point of the input u (which is usually determined by running power flow for every dispatch period), and $K \in \mathbb{R}^{n_u \times n}$ is the feedback controller gain matrix. The main idea in the given control law is to design a gain matrix K such that the power system converges back to its equilibrium and its transient stability can be improved after a large disturbance. To do that, one first needs to design the perturbed closed-loop dynamics whose states are then driven asymptotically to zero while minimizing the performance criterion associated with z . With that in mind, let us assume the power system model (2) is disturbed by an external uncertainty (which can be because of disturbance in load demand or renewable, etc.) and thus the system moves from its initial equilibrium status to a new status x' . Then, the perturbed closed loop dynamics and its associated performance criterion vector z can be written as follows:

$$E\Delta\dot{x} = (A - BK)\Delta x \quad (5a)$$

$$\Delta z = (C - DK)\Delta x \quad (5b)$$

where $\Delta x = x - x'$. The main objective of the matrix K is to drive Δx asymptotically to zero which in other words means that the power system (2) will converge to the new steady-state after large disturbance in system dynamics.

With little abuse of notation, for the sake of simplicity, from now on we drop the Δ notation from Δx , and Δz and simply use x and z instead, respectively. To that end, notice in (2) that conventional power plants, solar plants, wind power plants, and motor loads are dynamic systems while loads, non-unit buses, and other interconnections are considered static systems. Thus, the algebraic variables x_a in the power system model (the voltage and current phasors) are redundant states and thus can be eliminated while designing a feedback controller [2]. This can be done by finding an explicit equation for the algebraic variables using the algebraic constraint model and plugging it back into the system dynamical model as follows: Considering

$$A = \begin{bmatrix} A_{dd} & A_{da} \\ A_{ad} & A_{aa} \end{bmatrix}, x = \begin{bmatrix} x_d^T & x_a^T \end{bmatrix}^T, B = \begin{bmatrix} B_d^T & B_a^T \end{bmatrix}^T$$

one can write the DAE system (2) into the separate dynamic and algebraic model as follows:

$$\dot{x}_d = A_{dd}x_d + A_{da}x_a + B_{ud}u \quad (6a)$$

$$0 = A_{ad}x_d + A_{aa}x_a + B_{ua}u. \quad (6b)$$

Now assuming A_{aa} to be invertible which is a standard as-

sumption in power systems [2], [20]. Then an explicit equation for the algebraic variable x_a from (6b) can be written as:

$$x_a = A_{aa}^{-1}(-A_{ad}x_d - B_{ua}u). \quad (7)$$

Finally, based on this, the closed-loop model (5) can be rewritten in equivalent ODE format as follows:

$$\dot{x}_d(t) = (\tilde{A} - \tilde{B}K_d)x_d(t) \quad (8a)$$

$$z_1(t) = (\tilde{C} - \tilde{D}K_d)x_d(t) \quad (8b)$$

where $K_d \in \mathbb{R}^{n_u \times n_d}$ and the rest of the system matrices are defined as:

$$\tilde{A} = A_{dd} - A_{da}A_{aa}^{-1}A_{ad}, \quad \tilde{B} = B_d - A_{da}A_{aa}^{-1}B_a$$

$$\tilde{D} = B_d - C_aA_{aa}^{-1}B_a, \quad \tilde{C} = C_d - C_aA_{aa}^{-1}A_{ad}.$$

Notice that the computed gain K_d from ODE (8) can be plugged back into the complete interconnected power system as $K = [K_d \quad 0]$. Now, by just using the knowledge of the system matrices \tilde{A} , \tilde{B} , [21] has proposed an iterative algorithm to design a linear state feedback controller described as follows:

Theorem 1 ([21]): Let K_{d_0} be the initial stabilizing controller gain matrix such that $\tilde{A} - \tilde{B}K_{d_0}$ is Hurwitz and $P_k \in \mathbb{S}_{++}^{n_d \times n_d}$, $k = 0, 1, \dots$ be the unique solution to the following Lyapunov equation (which is linear in P_k)

$$A_c^T P_k + P_k A_c + Q + K_{d_k}^T R K_{d_k} - N K_{d_k} - (N K_{d_k})^T = 0 \quad (9)$$

with K_{d_k} being updated recursively at each iteration k using

$$K_{d_{k+1}} = R^{-1}(\tilde{B}^T P_k + N^T) \quad (10)$$

where $A_c = \tilde{A} - \tilde{B}K_{d_k}$, $Q = \tilde{C}^T \tilde{C}$, $R = \tilde{D}^T \tilde{D}$, and $N = \tilde{C}^T \tilde{D}$. Then, $\tilde{A} - \tilde{B}K_{d_k}$ is Hurwitz, and the solution of P_k as $k \rightarrow \infty$ converges to the optimal solution of the well-known algebraic Riccati equation (ARE) given as:

$$A_c^T P + P A_c - (P \tilde{B} + N) R^{-1} (\tilde{B}^T P + N^T) + Q = 0 \quad (11)$$

The above technique requires the complete knowledge of the system matrices \tilde{A} and \tilde{B} which is unrealistic as accurate models might not be available or be uncertain specifically in the case of power systems with high penetration of renewable energy resources. Then, to design a similar feedback gain K_d without the knowledge of the system matrices and just using the data (the system state trajectories which can be computed via PMU data, for example), we present the following problem statement:

Problem 1: Given the state trajectories $x_d(t)$ of the LODE system (8a) and its performance metric (8b), design feedback controller K_d without utilizing the knowledge of system matrices.

IV. RL-BASED FEEDBACK CONTROLLER DESIGN

Primarily based on the theory in [22], here we present model-free RL-based centralized feedback controller design. The presented approach computes the feedback controller gain matrix K_d purely based on the continuous interaction with the power system in real-time and does not require the information of system matrices. Thus, the presented technique has a major advantage over model-based feedback controller design as it can adaptively tune the parameters of its gain matrix based on the state information received.

To that end, we rewrite Eq. (8a) as

$$\dot{\mathbf{x}}_d(t) = \mathbf{A}_c \mathbf{x}_d(t) + \tilde{\mathbf{B}}(\mathbf{K}_{d_k} \mathbf{x}_d(t) + \mathbf{u}_d(t)) \quad (12)$$

where $\mathbf{u}_d = -\mathbf{K}_{d_k} \mathbf{x}_d$. The objective of the controller is to minimize the energy of the performance criterion vector \mathbf{z}_1 (8b) which can be written as

$$\begin{aligned} \int_0^\infty \mathbf{z}_1^\top \mathbf{z}_1 dt &= \int_0^\infty (\tilde{\mathbf{C}} \mathbf{x}_d - \tilde{\mathbf{D}} \mathbf{u}_d)^\top (\tilde{\mathbf{C}} \mathbf{x}_d - \tilde{\mathbf{D}} \mathbf{u}_d) dt \\ &= \int_0^\infty \mathbf{x}_d^\top \mathbf{Q} \mathbf{x}_d + \mathbf{u}_d^\top \mathbf{R} \mathbf{u}_d - 2 \mathbf{x}_d^\top \mathbf{N} \mathbf{u}_d dt \end{aligned}$$

with \mathbf{Q} , \mathbf{R} , and \mathbf{N} are given as in Theorem 1. Then, let us consider the optimal cost-to-go as $V(\mathbf{x}_d) = \mathbf{x}_d^\top \mathbf{P}_k \mathbf{x}_d$ where $\mathbf{P}_k \in \mathbb{S}_{++}^{n_d \times n_d}$. Notice that $V(\mathbf{x}_d)$ is a Lyapunov function and in the case of a model-based approach, \mathbf{P}_k can be computed via a recursive approach given in Theorem 1 which converges to the solution of the ARE.

With that in mind, taking the derivative of $V(\mathbf{x}_d)$ along the trajectories of (12), we then get

$$\dot{V} = \mathbf{x}_d^\top (\mathbf{A}_c^\top \mathbf{P}_k + \mathbf{P}_k \mathbf{A}_c) \mathbf{x}_d + 2(\mathbf{u}_d + \mathbf{K}_{d_k} \mathbf{x}_d)^\top \tilde{\mathbf{B}}^\top \mathbf{P}_k \mathbf{x}_d \quad (13)$$

From (9) let

$$\mathbf{X} = \mathbf{Q} + \mathbf{K}_{d_k}^\top \tilde{\mathbf{R}} \mathbf{K}_{d_k} - \mathbf{N} \mathbf{K}_{d_k} - (\mathbf{N} \mathbf{K}_{d_k})^\top. \quad (14)$$

Then, using (10) we can rewrite (13) as

$$\dot{V} = -\mathbf{x}_d^\top \mathbf{X} \mathbf{x}_d + 2(\mathbf{u}_d + \mathbf{K}_{d_k} \mathbf{x}_d)^\top (\mathbf{R} \mathbf{K}_{d_{k+1}} - \mathbf{N}^\top) \mathbf{x}_d. \quad (15)$$

With that in mind, taking integral on both sides of (15), we then have

$$\begin{aligned} \mathbf{x}_d^\top \mathbf{P}_k \mathbf{x}_d \Big|_t^{t+\delta t} &= \int_t^{t+\delta t} 2(\mathbf{u}_d + \mathbf{K}_{d_k} \mathbf{x}_d)^\top (\mathbf{R} \mathbf{K}_{d_{k+1}} - \mathbf{N}^\top) \mathbf{x}_d d\tau \\ &\quad - \int_t^{t+\delta t} \mathbf{x}_d^\top \mathbf{X} \mathbf{x}_d d\tau. \end{aligned} \quad (16)$$

Notice that Eq. (16) does not depend on any system matrices of the power system model. Then, if we can solve (16) for \mathbf{P}_k , $\mathbf{K}_{d_{k+1}}$ we have a solution to the posed Problem 1.

To that end, let us define matrices $\delta_{x_d x_d} \in \mathbb{R}^{s \times \frac{1}{2} n_d (n_d + 1)}$, $\Upsilon_{x_d x_d} \in \mathbb{R}^{s \times n_d^2}$, $\Upsilon_{x_d u} \in \mathbb{R}^{s \times n_d n_u}$, and $\hat{\mathbf{P}} \in \mathbb{R}^{\frac{1}{2} n_d (n_d + 1)}$ as follows:

$$\begin{aligned} \delta_{x_d x_d} &= [\bar{\mathbf{x}}_d|_{t_0}^{t_1}, \bar{\mathbf{x}}_d|_{t_1}^{t_2}, \bar{\mathbf{x}}_d|_{t_2}^{t_3}, \dots, \bar{\mathbf{x}}_d|_{t_{s-1}}^{t_s}]^\top \\ \Upsilon_{x_d x_d} &= \left[\int_{t_0}^{t_1} \mathbf{x}_d \otimes \mathbf{x}_d d\tau, \dots, \int_{t_{s-1}}^{t_s} \mathbf{x}_d \otimes \mathbf{x}_d d\tau \right]^\top \\ \Upsilon_{x_d u} &= \left[\int_{t_0}^{t_1} \mathbf{x}_d \otimes \mathbf{u}_d d\tau, \dots, \int_{t_{s-1}}^{t_s} \mathbf{x}_d \otimes \mathbf{u}_d d\tau \right]^\top \end{aligned} \quad (17)$$

$$\hat{\mathbf{P}} = [P_{11}, \dots, 2P_{1n_d}, P_{22}, \dots, 2P_{n_d-1n_d}, P_{n_d n_d}]^\top$$

where $s > 0$ are the total number of data samples collected, $\hat{\mathbf{P}}$ denotes the half vectorization of symmetric matrix \mathbf{P} with off-diagonal element taken as $2P_{ij}$, and vector $\bar{\mathbf{x}}_d \in \mathbb{R}^{\frac{1}{2} n_d (n_d + 1)}$ given as:

$$\bar{\mathbf{x}}_d = [\bar{x}_{d,1}^2, \dots, \bar{x}_{d,1} \bar{x}_{d,n_d}, \bar{x}_{d,2}^2, \dots, \bar{x}_{d,n_d-1} \bar{x}_{d,n_d}, \bar{x}_{d,n_d}^2]^\top$$

Then, (16) can equivalently be written as

$$\Xi_k \begin{bmatrix} \hat{\mathbf{P}}_k \\ \text{vec}(\mathbf{K}_{d_{k+1}}) \end{bmatrix} = \Psi_k \quad (18)$$

where $\Xi_k \in \mathbb{R}^{s \times (\frac{1}{2} n_d (n_d + 1) + n_d n_u)}$ and $\Psi_k \in \mathbb{R}^s$ are given as

$$\Xi_k = [\delta_{x_d x_d}, -2\Upsilon_{x_d x_d} (\mathbf{I}_{n_d} \otimes \mathbf{K}_{d_k}^\top \mathbf{R}) - 2\Upsilon_{x_d u} (\mathbf{I}_{n_d} \otimes \mathbf{R})]$$

$$\begin{aligned} \Psi_k &= -2\Upsilon_{x_d u} \text{vec}(\mathbf{N}) - 2\Upsilon_{x_d x_d} \text{vec}(\mathbf{K}_{d_k}^\top \mathbf{N}^\top) \\ &\quad - \Upsilon_{x_d x_d} \text{vec}(\mathbf{X}). \end{aligned}$$

To that end, Eq. (18) can efficiently be solved in least-squares sense for the variables $\hat{\mathbf{P}}_k$, $\text{vec}(\mathbf{K}_{d_{k+1}})$ as follows:

$$\begin{bmatrix} \hat{\mathbf{P}}_k \\ \text{vec}(\mathbf{K}_{d_{k+1}}) \end{bmatrix} = (\Xi_k^\top \Xi_k)^{-1} \Xi_k^\top \Psi_k. \quad (19)$$

If (19) has a solution, then the iterative model-based procedure to compute the feedback gain \mathbf{K}_d given by (9) and (10) in Theorem 1 can be replaced by model-free approach in (19). With that in mind, we now present a key assumption on the constructed data matrices $\Upsilon_{x_d x_d}$, $\Upsilon_{x_d u}$ which if gets satisfied guarantees the convergence of (19) to the solution of the corresponding ARE.

Assumption 1: The rank of matrix $[\Upsilon_{x_d x_d}, \Upsilon_{x_d u}]$ is equal to $(\frac{1}{2} n_d (n_d + 1) + n_d n_u)$.

Notice that Assumption 1 is satisfied by collecting sufficiently many samples and by adding exploration noise in the control inputs as explained in Section V.

Theorem 2: Given Assumption 1 is satisfied then the recursive solution of (19) as $k \rightarrow \infty$ converges to the optimal solution given by the corresponding ARE.

Proof: First, we need to prove that (19) has a unique solution which means we have to show that for

$$\Xi_k \begin{bmatrix} \hat{\mathbf{P}}_k \\ \text{vec}(\mathbf{K}_{d_{k+1}}) \end{bmatrix} = \mathbf{O} \quad (20)$$

$\hat{\mathbf{P}}_k$ and $\text{vec}(\mathbf{K}_{d_{k+1}})$ only has trivial solution, i.e., $\hat{\mathbf{P}}_k = \mathbf{O}$, $\text{vec}(\mathbf{K}_{d_{k+1}}) = \mathbf{O}$. This can be proved by contradiction. Let us assume $\mathbf{W} = [\mathbf{Y} \ \mathbf{Z}]^\top$ where $\mathbf{Y} \in \mathbb{R}^{n_d \times n_d}$ and $\mathbf{Z} \in \mathbb{R}^{n_u \times n_d}$, is a non-zero solution to (20). Now, for $\Upsilon_{\bar{x}_d} \in \mathbb{R}^{s \times \frac{1}{2} n_d (n_d + 1)}$ Eq. (18) can be expanded and rewritten as follows:

$$\begin{aligned} \Xi_k \begin{bmatrix} \hat{\mathbf{P}}_k \\ \text{vec}(\mathbf{K}_{d_{k+1}}) \end{bmatrix} &= -2\Upsilon_{x_d u} \text{vec}(\mathbf{N}) - 2\Upsilon_{\bar{x}_d} (\bar{\mathbf{N}}) - \Upsilon_{\bar{x}_d} (\hat{\mathbf{X}}) \\ &= [\Upsilon_{\bar{x}_d}, 2\Upsilon_{x_d u}] [(\hat{\mathbf{X}} - 2\bar{\mathbf{N}})^\top, \text{vec}(\mathbf{N})^\top]^\top \end{aligned}$$

or by (20) we have

$$[\Upsilon_{\bar{x}_d}, 2\Upsilon_{x_d u}] [(\hat{\mathbf{X}} - 2\bar{\mathbf{N}})^\top, \text{vec}(\mathbf{N})^\top]^\top = \mathbf{O}. \quad (21)$$

By Assumption 1 we know that $[\Upsilon_{\bar{x}_d}, 2\Upsilon_{x_d u}]$ has a full-column rank, which means that (21) only has trivial solution thus implying $\mathbf{N} = \mathbf{O}$ and $\mathbf{X} = \mathbf{O}$. Now from (9) with $\mathbf{X} = \mathbf{O}$ we have

$$\mathbf{A}_c^\top \mathbf{Y} + \mathbf{Y} \mathbf{A}_c = \mathbf{O}. \quad (22)$$

Since $\tilde{\mathbf{A}}_c$ is Hurwitz, then the only solution for (22) is $\mathbf{Y} = \mathbf{O}$. Then, from (10) with $\mathbf{Y} = \mathbf{O}$ and $\mathbf{N} = \mathbf{O}$ we have $\mathbf{Z} = \mathbf{O}$, thus $\mathbf{W} = \mathbf{O}$ which contradicts with our assumption that $\mathbf{W} \neq \mathbf{O}$. Thus, a unique solution exists for Eq. (19).

With that in mind, since $\hat{\mathbf{P}}_k$ and $\text{vec}(\mathbf{K}_{d_{k+1}})$ can uniquely be determined, then the policy iteration in (19) is equivalent to (9) and (10) and thus by Theorem 1 the convergence to the ARE solution is guaranteed. This completes the proof. ■

Computing controller gain matrix \mathbf{K}_d by solving (19) gives us a model-free approach to designing an optimal feedback controller. The presented approach does not require any information regarding the power system model and thus is an

Algorithm 1: RL-WADC for Interconnected Power Systems

- 1 **Save data:** Simulate power system model (8) with control input $\mathbf{u}(t) = \mathbf{u}_0 + \mathbf{e}(t)$ and construct the data matrices given (17) until Assumption 1 is satisfied.
 - 2 **Compute \mathbf{K}_d iteratively:** Starting with initial stabilizing controller \mathbf{K}_{d_0} compute \mathbf{K}_{d_k} recursively for $k = 0, 1, \dots$ via solving the following equation till $|\hat{\mathbf{P}}_k - \hat{\mathbf{P}}_{k-1}| < \epsilon$ where ϵ is a user-defined threshold
$$\begin{bmatrix} \hat{\mathbf{P}}_k \\ \text{vec}(\mathbf{K}_{d_{k+1}}) \end{bmatrix} = (\Xi_k^\top \Xi_k)^{-1} \Xi_k^\top \Psi_k.$$
 - 3 **Apply the controller:** Remove the exploration noise $\mathbf{e}(t)$ from the control input, design \mathbf{K} as $\mathbf{K} = [\mathbf{K}_d \ \mathbf{O}]$, and apply the following control policy to the complete NDAE power system model
$$\mathbf{u}_C(t) = \mathbf{u}_0^k - \mathbf{K} (\mathbf{x}(t) - \mathbf{x}^k).$$
-

adaptive approach that can finely tune its controller gain matrix depending on the received data. This is especially useful for the power grid with uncertain renewable resources and dynamic loads. In the following section, we assess the performance of the presented RL-WADC through numerical case studies on an interconnected power system model with detailed renewable dynamics and various composite loads.

V. CASE STUDIES

The presented RL-WADC has been tested on a modified IEEE 9-bus test network. This test system consists of one conventional steam power plant at Bus 1, a solar farm at Bus 2, a wind power plant at Bus 3, constant power and constant impedance loads at Buses 5 and 6, and a motor load at Bus 8. Further details about the parameters of solar/wind farms, conventional power plants, and dynamics of the test system used in this study can be found in [1], [19], [20].

All the simulations are carried out on a personal computer with 64GB RAM and an Intel i9 – 11900K processor. The power system NDAE dynamics have been modeled in MATLAB R2021a and are simulated using `ode15s` MATLAB DAEs solver. The initial conditions for the `ode15s` are computed using power flow studies in MATPOWER through function `runpf`. The volt-ampere base for the power system is chosen to be $S_b = 100\text{MVA}$ while the frequency base is selected to be $\omega_b = 120\pi\text{rad/s}$. It is worth mentioning here that in the above test system, to meet the overall load demand in the steady state, on average 63% of the power is generated by the renewables (both solar and wind power plants), hence it can be considered as a renewable heavy power grid.

To design the RL-WADC based feedback controller gain matrix \mathbf{K}_d the main requirement is to compute the data matrices $\delta_{x_d x_d}$, $\Upsilon_{x_d x_d}$, $\Upsilon_{x_d u}$ as given in (17). To do that, the power system model is simulated with exploration noise in the initial control policy and the data matrices are saved until Assumption 1 is satisfied. Notice that, since the power system model is already stable, the initial control gain \mathbf{K}_{d_0} is

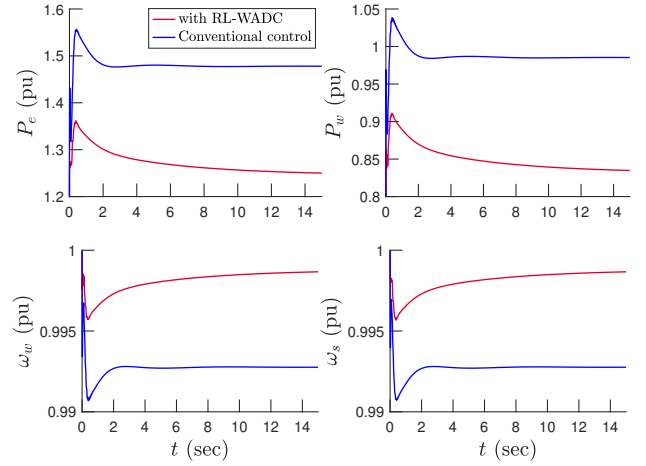


Fig. 1. Performance under abrupt load disturbance: Active power injected and relative speed of solar and wind power plant, respectively.

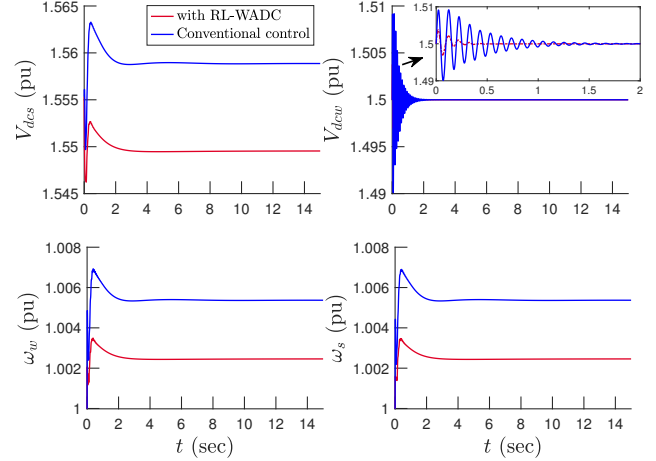


Fig. 2. Performance under renewable and load disturbance: DC link voltage and relative speed of both solar and wind power plant.

then chosen to be zero. The selection of the type of exploration noise is not trivial in RL-based and other machine learning-based techniques and various types of noises have been used in the literature such as sum of sinusoids, exponentially decreasing noise, and random bounded noise and vice versa. In this work, we have used a sum of sinusoids with different frequencies as exploration noise in the initial control policy as given in [22]. The overall implementation of the presented RL-WADC is summarized in Algorithm 1. Notice that in the case of parametric uncertainty the Algorithm 1 can simply be repeated to adjust gain matrix \mathbf{K} accordingly.

A. Performance under transient conditions

Here we assess the performance of the presented RL-WADC under transient conditions created by various disturbances in renewables and load demand. To that end, we initially test the performance of RL-WADC by adding, right after $t > 0$, a random step increasing in load demand at Bus 6. This can be a result of a sudden generator disconnection or a transmission line trip, hence increasing the overall load demand of the power network. The system dynamic response with just system primary controllers and with RL-WADC acting on top of them is shown in Fig. 1. Notice that the

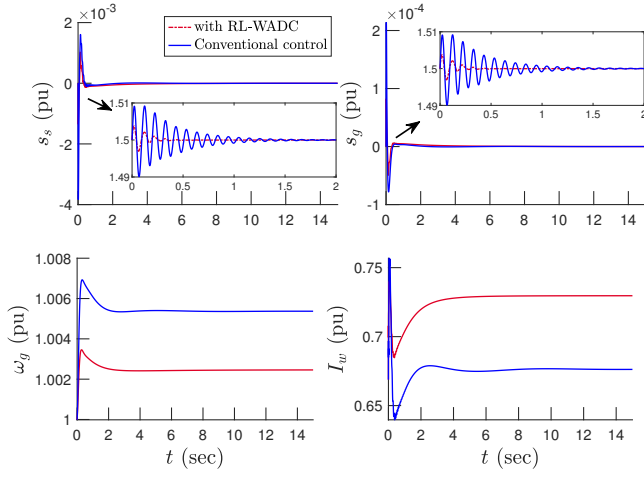


Fig. 3. Performance under renewable and load disturbance: solar plant slip (top left), generator slip (top right) and frequency (bottom left), and wind power plant current (bottom right).

primary/conventional controllers in the test power system include an automatic voltage regulator (AVR), governor, and power system stabilizer (PSS) for conventional power plants [20] and proportional-integral (PI) type controllers for both solar and wind power plants [19], [20]. The proposed RL-WADC acts on top of these primary controllers and sends additional control signals via input u . The overall objective of the proposed RL-WADC is to improve system transient stability by adding damping to the system oscillation and bringing the system back close to its nominal/equilibrium value after a large disturbance.

To that end, we can see from Fig. 1 that without the proposed controller the power output from both solar and wind power plants have oscillations while with RL-WADC there is significant damping. Similarly, for the frequencies of both solar and wind power plants, we can see that without RL-WADC after disturbance their frequencies dip close to 0.99 (pu) while with the proposed controller the overall dip is around 0.994 (pu), thus improving system frequency nadir.

To further advocate the performance of the presented RL-based damping controller, here we add further disturbance this time by decreasing load demand by 20% and solar irradiance by 10% on the solar power plant. The results are presented in Figs. 2 and 3. We can see from Fig. 2 that with the proposed controller after disturbance there is significant damping in the slip of synchronous generator, wind, and solar power plant showing significant improvements in LFOs and ULFOs. Similarly, from Fig. 3 we can see that there is a significant damping in the system oscillations thus the overall transient stability after the large disturbance has been improved.

REFERENCES

- [1] M. Nadeem, M. Bahavarnia, and A. F. Taha, "On wide-area control of solar-integrated dae models of power grids," in *American Control Conference (ACC)*. IEEE, 2023, pp. 4495–4500.
- [2] T. Sadamoto, A. Chakraborty, T. Ishizaki, and J.-i. Imura, "Dynamic modeling, stability, and control of power systems with distributed energy resources: Handling faults using two control methods in tandem," *IEEE Control Systems Magazine*, vol. 39, no. 2, pp. 34–65, 2019.
- [3] X. Chen, G. Qu, Y. Tang, S. Low, and N. Li, "Reinforcement learning for selective key applications in power systems: Recent advances and future challenges," *IEEE Transactions on Smart Grid*, vol. 13, no. 4, pp. 2935–2958, 2022.
- [4] A. Chakraborty, "Wide-area control of power systems: Employing data-driven, hierarchical reinforcement learning," *IEEE Electrification Magazine*, vol. 9, no. 1, pp. 45–52, 2021.
- [5] A. F. Dizche, A. Chakraborty, and A. Duel-Hallen, "Sparse wide-area control of power systems using data-driven reinforcement learning," in *2019 American Control Conference (ACC)*, 2019, pp. 2867–2872.
- [6] T. Yu, B. Zhou, K. W. Chan, L. Chen, and B. Yang, "Stochastic optimal relaxed automatic generation control in non-markov environment based on multi-step $q(\lambda)$ learning," *IEEE Transactions on Power Systems*, vol. 26, no. 3, pp. 1272–1282, 2011.
- [7] W. Cui, Y. Jiang, and B. Zhang, "Reinforcement learning for optimal primary frequency control: A lyapunov approach," *IEEE Transactions on Power Systems*, vol. 38, no. 2, pp. 1676–1688, 2022.
- [8] J. Feng, W. Cui, J. Cortés, and Y. Shi, "Online event-triggered switching for frequency control in power grids with variable inertia," *arXiv preprint arXiv:2408.15436*, 2024.
- [9] S. Wang, J. Duan, D. Shi, C. Xu, H. Li, R. Diao, and Z. Wang, "A data-driven multi-agent autonomous voltage control framework using deep reinforcement learning," *IEEE Transactions on Power Systems*, vol. 35, no. 6, pp. 4644–4654, 2020.
- [10] Z. Yan and Y. Xu, "A multi-agent deep reinforcement learning method for cooperative load frequency control of a multi-area power system," *IEEE Transactions on Power Systems*, vol. 35, no. 6, pp. 4599–4608, 2020.
- [11] H. Wang, Z. Lei, X. Zhang, J. Peng, and H. Jiang, "Multiobjective reinforcement learning-based intelligent approach for optimization of activation rules in automatic generation control," *IEEE access*, vol. 7, pp. 17 480–17 492, 2019.
- [12] V. P. Singh, N. Kishor, and P. Samuel, "Distributed multi-agent system-based load frequency control for multi-area power system in smart grid," *IEEE Transactions on Industrial Electronics*, vol. 64, no. 6, pp. 5151–5160, 2017.
- [13] Y. Shi, G. Qu, S. Low, A. Anandkumar, and A. Wierman, "Stability constrained reinforcement learning for real-time voltage control," in *2022 American Control Conference (ACC)*. IEEE, 2022, pp. 2715–2721.
- [14] L. Yin, C. Zhang, Y. Wang, F. Gao, J. Yu, and L. Cheng, "Emotional deep learning programming controller for automatic voltage control of power systems," *IEEE Access*, vol. 9, pp. 31 880–31 891, 2021.
- [15] M. Glavic, "(deep) reinforcement learning for electric power system control and related problems: A short review and perspectives," *Annual Reviews in Control*, vol. 48, pp. 22–35, 2019.
- [16] S. Mukherjee, A. Chakraborty, H. Bai, A. Darvishi, and B. Fardanesh, "Scalable designs for reinforcement learning-based wide-area damping control," *IEEE Transactions on Smart Grid*, vol. 12, no. 3, pp. 2389–2401, 2021.
- [17] D. Vrabie, O. Pastravanu, M. Abu-Khalaf, and F. L. Lewis, "Adaptive optimal control for continuous-time linear systems based on policy iteration," *Automatica*, vol. 45, no. 2, pp. 477–484, 2009.
- [18] Y. Jiang and Z.-P. Jiang, "Robust adaptive dynamic programming for large-scale systems with an application to multimachine power systems," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 59, no. 10, pp. 693–697, 2012.
- [19] H. N. V. Pico and V. Gevorgian, "Blackstart capability and survivability of wind turbines with fully rated converters," *IEEE Transactions on Energy Conversion*, vol. 37, no. 4, pp. 2482–2497, 2022.
- [20] M. Nadeem, M. Bahavarnia, and A. F. Taha, "Robust feedback control of power systems with solar plants and composite loads," *IEEE Transactions on Power Systems*, vol. 39, no. 3, pp. 4949–4962, 2024.
- [21] W. Arnold and A. Laub, "Generalized eigenproblem algorithms and software for algebraic riccati equations," *Proceedings of the IEEE*, vol. 72, no. 12, pp. 1746–1754, 1984.
- [22] Y. Jiang and Z.-P. Jiang, "Computational adaptive optimal control for continuous-time linear systems with completely unknown dynamics," *Automatica*, vol. 48, no. 10, pp. 2699–2704, 2012.