

Highlights

Reinforcement Learning Optimizes Power Dispatch in Decentralized Power Grid

Yongsun Lee, Hoyun Choi, Laurent Pagnier, Cook Hyun Kim, Jongshin Lee, Bukyoung Jhun, Heetae Kim, Jürgen Kurths, B. Kahng

- The advent of renewable sources has introduced instability into modern power grids.
- Maintaining the frequency of grid is imperative, even in the power plant accidents.
- It is challenging to decide how much unbalanced power each bus should compensate.
- We address the issue through deep reinforcement learning with graph neural network.

Reinforcement Learning Optimizes Power Dispatch in Decentralized Power Grid

Yongsun Lee^a, Hoyun Choi^{a,b}, Laurent Pagnier^c, Cook Hyun Kim^b, Jongshin Lee^b, Bukyoung Jhun^a, Heetae Kim^d, Jürgen Kurths^{b,e,f}, B. Kahng^{b,*}

^a*CTP and Department of Physics and Astronomy, Seoul National University, Seoul, 08826, Korea*

^b*CCSS, KI for Grid Modernization, Korea Institute of Energy Technology, Naju, 58217, Jeonnam, Korea*

^c*Department of Mathematics, The University of Arizona, Tucson, 85721, Arizona, USA*

^d*KI for Grid Modernization, Korea Institute of Energy Technology, Naju, 58217, Jeonnam, Korea*

^e*Potsdam Institute for Climate Impact Research, Telegraphenberg, D-14415, Potsdam, Germany*

^f*Institute of Physics, Humboldt University Berlin, Berlin, D-12489, Germany*

Abstract

Effective frequency control in power grids has become increasingly important with the increasing demand for renewable energy sources. Here, we propose a novel strategy for resolving this challenge using graph convolutional proximal policy optimization (GC-PPO). The GC-PPO method can optimally determine how much power individual buses dispatch to reduce frequency fluctuations across a power grid. We demonstrate its efficacy in controlling disturbances by applying the GC-PPO to the power grid of the UK. The performance of GC-PPO is outstanding compared to the classical methods. This result highlights the promising role of GC-PPO in enhancing the stability and reliability of power systems by switching lines or decentralizing grid topology.

Keywords: Reinforcement learning, Graph neural network, Decentralized power grid, Power dispatch

1. Introduction

As global warming accelerates, our energy systems must transition rapidly and substantially to those that utilize renewable energy sources, such as biomass, hydroelectric power, wind, and solar energy. Electric power grids, a fundamental component of this transition, are rapidly adopting renewable energy sources. However, as the proportion of renewable energy in the grid increases, grid systems will become more susceptible to instability [1, 2, 3, 4, 5, 6].

Challenges abound in managing such hybrid power systems. Solar and wind energy generation depends heavily on daily and seasonal weather patterns. Furthermore, the electric frequency can fluctuate within seconds, owing to variations in sunlight and wind intensity caused by clouds and storms. Consequently, regulating electric currents from renewable sources is becoming increasingly vital for stabilizing these frequencies [7, 8, 9, 10, 11]. Additionally, renewable energy power plants are often far from consumer loca-

tions, necessitating long-distance transmission lines and posing a cascading failure risk [12, 13].

Therefore, optimizing decentralized hybrid power grids has emerged as a critical consideration. Power grids transmit two types of alternating current (AC): inertial AC, derived from fossil fuel combustion or nuclear fission, and inertia-free AC, generated by renewable sources and connected via power electronic inertia-free inverters [14, 15, 16, 17]. External disturbances such as transmission line disconnections or regional overloads can disrupt the grid, with inertia-free AC unable to recover spontaneously, potentially leading to widespread frequency desynchronization and grid failure [18, 19, 20, 21, 22].

Considerable efforts have been made to reduce these dynamic disturbances and avoid large-scale power grid blackouts. Several methods have been proposed and implemented, such as controlling the time-dependent feedback (e.g., fast frequency responses [23]), increasing the global inertia by connecting turbines without generators [24, 25] and switching off uncontrollable generators [26]. A traditional method of recovering voltage or frequency is to add extra power saved in

*corresponding author

Email address: bkahng@kentech.ac.kr (B. Kahng)

electric storage systems [27, 28] or request plants to produce extra power calculated using the optimal reactive power dispatch algorithm [29, 30]. These methods have been effectively applied for a long time because the power grid is centralized, and the number of plants is small. However, in power grids composed of many small solar and wind plants, such strategies may not be optimal owing to their slow response and limited adaptability to rapid changes in power generation [31, 32]. Thus, finding an optimal method to recover instantaneously or dispatch power to stabilize the grid system is a significant challenge.

In this context, reinforcement learning (RL) has emerged as a promising approach for devising optimal dynamic strategies. RL is often used to determine optimal dynamic pathways in various fields, ranging from the game of Go to autonomous driving [33, 34, 35, 36, 37], and provides efficient algorithms for diverse phenomena. Accordingly, we propose a novel power dispatch strategy using the RL approach, specifically the Graph Convolutional Proximal Policy Optimization (GC-PPO) algorithm. Using this method, we obtain information on the amount of extra power produced by each generator to minimize frequency fluctuations across the grid, offering quick adaptability to diverse grid configurations. The extra power supply of each generator is heterogeneous and depends on the topology of the power grid.

2. Results

2.1. Swing equation of oscillators in the power grid

We consider a power grid comprising N_g generators and N_c consumers. Therefore, the power grid comprises $N = N_g + N_c$ buses. Their phases and frequencies are denoted as $(\theta_i, \dot{\theta}_i)$, with index $i = 1, \dots, N$. The buses are connected to other buses via transmission lines and are treated as nodes in the graph representation. An oscillator at bus i rotates according to the swing equation (also called the second-order Kuramoto model) [38, 39]

$$m_i \ddot{\theta}_i + \gamma_i \dot{\theta}_i = P_i + \sum_{j=1}^N K_{ij} \sin(\theta_j - \theta_i), \quad (1)$$

where m_i is the angular momentum (or called inertia), $\ddot{\theta}_i$ is the angular acceleration, P_i is the amount of power generation ($P_i > 0$) or consumption ($P_i < 0$), γ_i is the damping coefficient of oscillator i , and K_{ij} is the coupling constant (or called line susceptance)

between connected buses i and j , and is given as $K_{ij} = |V_i||V_j|/x_{ij}$, where V_i and V_j are the voltages of buses i and j , and x_{ij} is the reactance of the transmission line. If buses i and j are not connected, then $K_{ij} = 0$.

2.2. Network models

We test our method on different networks to validate it. Here, we use both synthetic SHK networks [40] that have the advantage of being easily generated and a realistic model of the UK grid, which captures more faithfully the details of real-world power systems. The SHK model is designed to reflect the features of real-world grids from various aspects. The topology is governed by a trade-off between a tree-like structure to minimize costs and ensuring redundant backup lines for emergencies.

Physical variables in the SHK grid are taken as $m_i = m = 1$, $\gamma_i = \gamma = 0.5$, and power is chosen as $P_i = P = 1$ (generators) and $P_i = P = -1$ (consumers). The proportions of generators and consumers are even, whereas their positions on the graph are randomly assigned. $K_{ij} = K = 4$.

The real-world power grid of the UK comprises 235 buses, where generators control their power generation, and consumers regulate their consumption, for example, with fast frequency response. The power grid is reduced to 54 buses comprising $N_g = 25$ renormalized generators and $N_c = 29$ renormalized consumers using the Kron reduction method [41] for computation. Since these reduced buses consist of multiple generators and consumers, their power $\{P_i\}$ can be controlled. The topology of the reduced UK grid is illustrated in Fig. 1(a), where the red diamond buses represent generators ($P > 0$) and the blue circles represent consumers ($P < 0$). The power, mass, and damping coefficient of each bus are considered as their physical values (Figs. 1 (f)–(h)). The UK grid comprises buses with different inertia, including renewable sources, and transmission lines with different coupling constants. As a result, it has more complex dynamics than the synthetic power grid.

2.3. Bus-based power dispatch

To determine the steady synchronous state of a given power grid, we first randomly select θ_i in the range $(-\pi, \pi]$ and $\dot{\theta}_i = 0$ for each $i = 1, \dots, N$. Then, the nodes' $\{\theta_i, \dot{\theta}_i\}$ values are updated following the swing equation (1) until they reach a steady state. The power-dispatch process begins with the initial

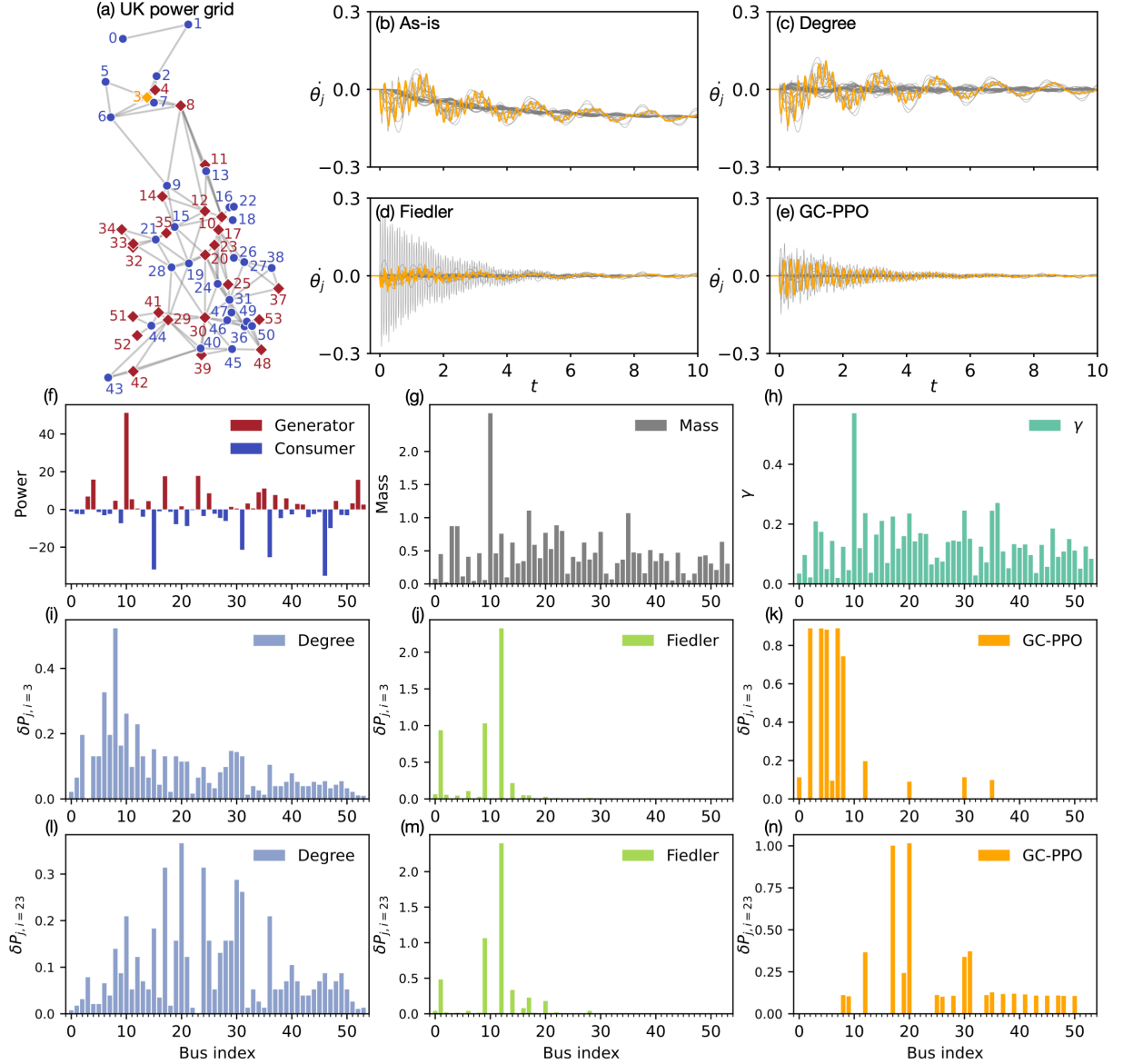


Figure 1: (a) Power grid of the high-voltage transmission lines in the UK coarse-grained by the Kron reduction method. Buses are composed of generators (Red diamond) with $P_i > 0$ and consumers (blue circle) with $P_i < 0$. (b) The frequency relaxation pattern of each bus after bus 3 (yellow diamond) is perturbed, and no power dispatch is applied. The frequencies for the other buses are drawn in gray. (c)–(e) Fluctuation relaxation pattern when the amount of power dispatch is determined by three different protocols: (c) Degree, (d) Fiedler, and (e) GC-PPO. Here are plots of the (f) power, (g) mass (inertia), and (h) damping coefficient versus bus indices of the UK grid. (i)–(k) Amount of power dispatch generated from bus j when bus $i = 3$ is perturbed, i.e., $\delta P_{j,i=3}$ for the three different protocols. (l)–(n) Similar to plots (i)–(k), but when bus 23, which is located at the center of the grid, is perturbed.

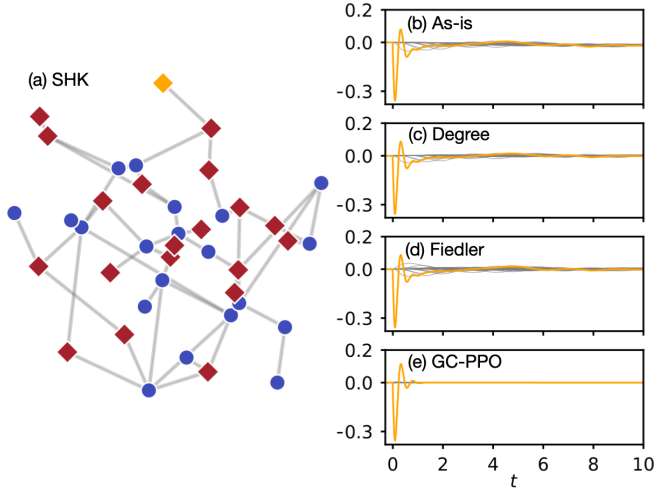


Figure 2: (a) Topology of the synthetic SHK grid, composed of generators (Red diamond) with $P_i = 1$ and consumers (blue circle) with $P_i = -1$. (b) Frequency evolution of the perturbed bus (Yellow line, denoted by Yellow diamond in (a)) and the other buses (gray lines) without power dispatch. (c)–(e) Fluctuation relaxation pattern when power dispatch is performed following three protocols: (c) Degree, (d) Fiedler, and (e) GC-PPO.

configuration of the steady state $\{\theta_i, \dot{\theta}_i\}$. Next, we consider two situations where a stable power grid is perturbed: (i) when a generator malfunctions or a consumer overuses power. For failure in generator i , the reduced power generation δP_i is set as, for example, $0.3P_i$. In addition to the change in power, the inertia of generator m_i is changed to $m_0 + 0.7(m_i - m_0)$. m_0 is set to 0.02 for the SHK grid and 0.025974 for the UK grid, which is the smallest inertia in the system. (ii) In the case of an overload, only the power consumption of the consumer i increases by $\delta P_i = 0.3|P_i|$, but the inertia m_i remains unchanged. Finally, to maintain the power balance $\sum_k P_k = 0$, each remaining bus P_j ($j \neq i$) must generate additional power denoted by δP_{ji} , which satisfies $\sum_{j \in \mathcal{B} \setminus \{i\}} \delta P_{ji} = \delta P_i$. Here, \mathcal{B} is defined as the set of all buses. δP_{ji} can be determined according to various protocols. The details are described in METHOD section 4.

Fig. 1(b) shows the evolution of the frequency of each bus when a perturbation is applied at $t = 0$ to generator 3, which is marked by a yellow diamond in Fig. 1(a), of the UK grid without power dispatch (As-is). Bus 3 exhibits the most severe fluctuation when no power dispatch is performed. The yellow line indicates the frequency of the perturbed bus 3, whereas the gray lines are the frequencies of the other buses. Figs. 1(c)–(e) show the frequency patterns of bus 3

and the other buses after power dispatch is applied following the Degree, Fiedler, and GC-PPO protocols, respectively. The value of $\delta P_{j,i=3}$ is given according to each protocol, as shown in Figs. 1(i)–(k). Because all buses except the perturbed one participate in the power dispatch, only the disturbed value $\delta P_{i,i}$ is zero. Although all protocols reduced fluctuations, the pattern of fluctuation relaxation suggests that the GC-PPO protocol achieves the best improvement among the protocols we tested.

Figs. 1(l)–(n) show plots similar to those shown in Figs. 1(i)–(k), when generator 23, located almost at the center of the grid, is perturbed. Each protocol responds differently, depending on perturbation δP_i . As shown in both Figs. 1(k) and (n), GC-PPO dispatches more power to generators in the neighborhood of the perturbed bus i , although it does not impose any constraints on the distance. In contrast, the Degree protocol redistributes power relatively evenly among the other generators. The Fiedler protocol focuses on a few buses (1, 9, 12) regardless of the perturbed bus.

Figs. 2 are corresponding plots to Figs. 1(a)–(e), but the network is adopted from the synthetic SHK power grid. Since the SHK grid is more homogeneous than the UK grid in topology, masses, and links aspects, the fluctuations of the buses are not so large compared to those of the UK grid. Nevertheless, when GC-PPO protocol is employed, the system quickly recovers the stable steady state.

2.4. Fluctuation measure

If power dispatch is successfully implemented, all buses will be synchronized with the previous frequency, and the grid will be stabilized. Otherwise, buses would not be in a synchronized state, and a global blackout could occur. Therefore, the stability of the power grid must be measured accurately after dispatch [42, 43].

Here, we introduce Ξ_i , a fluctuation measure of the grid after the power is dispatched owing to a perturbation on bus i .

$$\Xi_i \equiv \frac{1}{T} \int_0^T dt \left[\frac{1}{\sum_k m_k} \sum_k m_k \dot{\theta}_k^2(t) - \left(\frac{1}{\sum_k m_k} \sum_k m_k \dot{\theta}_k(t) \right)^2 \right] \quad (2)$$

The fluctuation measure Ξ_i is the weighted variance of the frequencies $\dot{\theta}_k(t)$ over all buses, including the generators and consumers. The weight is taken as inertia m_k [44]. The frequencies are monitored for T seconds after the initial perturbation at $t = 0$.

Power grid	Fluctuation measure Ξ			
	SHK		UK	
	Generator	Consumer	Generator	Consumer
As-is	5.25×10^{-3}	4.92×10^{-3}	1.71×10^{-1}	1.48×10^{-1}
Uniform	4.29×10^{-3}	4.12×10^{-3}	3.89×10^{-2}	2.92×10^{-2}
Degree	4.32×10^{-3}	4.19×10^{-3}	3.85×10^{-2}	2.92×10^{-2}
BC	4.43×10^{-3}	4.31×10^{-3}	3.73×10^{-2}	3.27×10^{-2}
Clustering	6.64×10^{-3}	5.74×10^{-3}	4.30×10^{-2}	2.87×10^{-2}
Fiedler	6.26×10^{-3}	6.19×10^{-3}	4.58×10^{-2}	5.98×10^{-2}
GC-PPO	2.54×10^{-3}	2.34×10^{-3}	1.13×10^{-2}	1.07×10^{-2}

Table 1: Total fluctuation measures Ξ over all buses for different protocols for the synthetic network and real-world UK power grid. Note that a smaller value of Ξ indicates a better protocol. The GC-PPO protocol is more efficient for the more heterogeneous UK power grid.

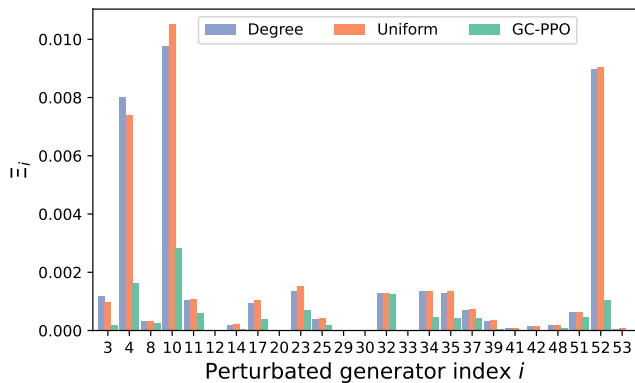


Figure 3: Fluctuation measure Ξ_i versus generator index i in the UK grid for three protocols: Degree, Uniform, and GC-PPO. Degree and Uniform protocols perform poorly on three generators (4, 10, and 52), while GC-PPO effectively moderates the fluctuations.

As the generator or consumer that causes the unexpected perturbation is not recognized, we introduce the average Ξ_i over all possible perturbations $\Xi \equiv \sum_{i \in \mathcal{G}} \Xi_i / N_g$ for the generators and the $\sum_{i \in \mathcal{C}} \Xi_i / N_c$ for the consumers. Naturally, the smaller Ξ is, the more stable the power grid and the better the protocol performance.

Fig. 3 shows the fluctuation measure Ξ_i versus the perturbed generator index i for the three protocols: Degree, Uniform, and GC-PPO. For example, for $i = 3$, Ξ_i of the GC-PPO is obtained based on the fluctuating data $\dot{\theta}(t)$ shown in Fig. 1 (e). While the Degree and Uniform protocols result in a similar pattern based on the peaks of $i = 4, 10$, and 52 , the GC-PPO protocol produces significantly smaller fluctuations. Generator 4 is in a sparsely connected region, and generator 52 is located on a leaf, making them topologically vulnerable nodes. Bus 10 is a gen-

erator with large P_i, m_i, γ_i as shown in Figs. 1 (f)–(h) and therefore, the perturbation caused by $0.3P_i$ significantly affects the stability of the entire system. Overall, the new GC-PPO protocol is more effective than the other protocols in reducing the instability of the UK power grid.

The Ξ values are listed in Table. 1, for six protocols to compare their capabilities with the raw values obtained without any protocol. In the case of the SHK grid, the fluctuation is small even when power dispatch is not performed (As-is). The physical quantities such as power, mass, damping coefficient, and coupling constants are rather homogeneous. Nevertheless, GC-PPO yields Ξ values several times smaller than the other protocols. GC-PPO demonstrates a considerably more dramatic performance improvement for the UK grid with heterogeneous physical quantities than for the synthetic grid. Whereas the other protocols exhibit Ξ improvements of approximately 10 times over no power dispatch (As-is), GC-PPO reduces Ξ by more than 100 times.

2.5. Training of GC-PPO

In contrast to the other protocols, the new GC-PPO protocol requires training. During training, the GC-PPO learns by experiencing various power dispatches and refining the method it uses to compute $\delta P_{j,i}$. While the detailed training method is described in Section 4.B, observing the progress of the GC-PPO performance would be beneficial. Fig. 4 shows a decrease in Ξ as the training of GC-PPO progresses. The dotted lines represent the Ξ values for other protocols in Table. 1. Because the other protocols do not require training, their Ξ s values remain constant regardless of the training episodes of the GC-PPO.

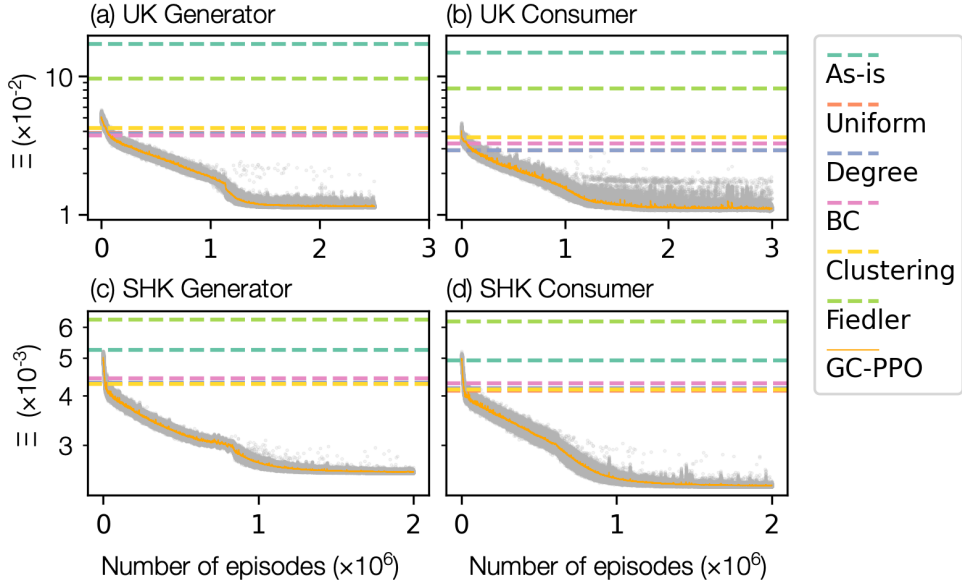


Figure 4: Performance of GC-PPO versus the number of training episodes in different environments for the different power grids and bus types. The solid yellow curves indicate the average performance of GC-PPO, whereas the gray dots represent the fluctuations due to the stochastic training process. The dashed lines are the performances of the other topology-based protocols for reference.

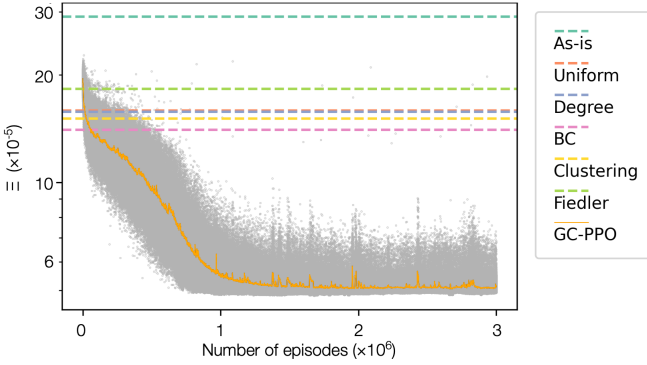


Figure 5: Training of the GC-PPO when multiple consumers simultaneously overuse power. Solid and dashed lines are equivalent to Fig. 4.

Across all power grids and perturbations, the GC-PPO clearly outperforms the other protocols, with a quick drop of Ξ in the early stage of training. It then shows a gradual improvement in performance until it reaches 10^6 training episodes. Finally, Ξ reaches a plateau where the improvement is negligible, indicating that the GC-PPO has been sufficiently trained.

Power shortages may occur when heaters or coolers are used simultaneously in a town or region be-

cause of significant weather changes. We envision a situation where six consumers in the northern part of the UK grid with bus indices 0, 1, 2, 5, 6, and 7 overuse power simultaneously. Fig. 5 shows the training curve for GC-PPO performance under this perturbation. Similar to Fig. 4, this figure reveals that GC-PPO outperforms the baseline protocols after a few training episodes, and its performance gradually improves over approximately 10^6 training cycles.

3. Discussion

We developed a novel optimal power-dispatch method using the GC-PPO algorithm classified into reinforcement learning. This method determines how much power *each bus* should dispatch a power to compensate for the lost power to stabilize the system. This method is compared to the traditional method of supplying the lost power from a few power plants. This paradigm shift in the power dispatching method has been a timely need as the power grid has become decentralized due to the proliferation of small-scale solar and wind power plants. The innovative algorithm of GC-PPO introduced here outperforms the classical

method.

The GC-PPO algorithm optimizes a variation measure Ξ , which is the weighted variance of the frequency over all buses. The weights are the inertia m_k of each bus k . This weighted strategy is more effective than the unweighted strategy in obtaining additional power $\delta P_{j,i}$, particularly for a heterogeneous power grid, meaning generators with larger inertia can be asked to produce more power to recover the system quickly. On the other hand, for the SHK model, we assume uniform inertia for all buses, i.e., $m_i = m = 1$. Then, no weights are needed to measure fluctuations.

Notably, GC-PPO is not limited to the specific situation we have covered; to implement the algorithm in other grid systems, we first pre-train GC-PPO on an ensemble of small power grids and then fine-tune it for a given large power grid [45, 46, 47]. Perturbations are not necessarily limited to a single bus but also occur when multiple consumers overuse simultaneously. Furthermore, GC-PPOs can also be proactive in grid stability if they are informed about possible faults [48].

Changing the grid's topology by switching lines is an alternative to maintain the stability of the power grid. The next challenge will be to modify the GC-PPO algorithm to solve the switching problem.

4. Methods

Here, we describe the determination of the power dispatch δP_{ji} owing to perturbation δP_i for each protocol listed in Table. 1. Because the power balance relation $\sum_{j \in \mathcal{G} \setminus \{i\}} \delta P_{ji} = \delta P_i$ should be satisfied, we choose δP_{ji} as proportional to δP_i .

$$\delta P_{ji} = \frac{q_{ji}}{\sum_{k \in \mathcal{G} \setminus \{i\}} q_{ki}} \delta P_i. \quad (3)$$

Thus, each protocol determines q_{ji} for generator j .

4.1. Heuristics methods

We propose five protocols based on the structural characteristics of the power grid: uniform, degree, betweenness centrality (BC) [49, 50], clustering coefficients [51], and spectral properties (Fiedler mode) [52], which are important for assessing electric power grid vulnerability [53]. These approaches possess different properties. We use the shortest-path distance d_{ji} between the perturbed bus i and compensating bus j . Because the flux per unit length decays by $1/d$ with

distance d from a source in two dimensions, we choose q_{ji} as inversely proportional to d_{ji} .

One of the simplest power dispatch protocols is the uniform protocol. All non-perturbed generators are treated equally, thus $q_{ji} = 1/d_{ji}$. Note that this protocol uses only distance information from the perturbed bus.

The degree of bus i (k_i) is defined as the number of buses connected to it via transmission lines. A bus with a large degree can exhibit increased stability when it is well-balanced with neighbors. Therefore, we defined a degree protocol $q_{ji} = k_j/d_{ji}$.

The BC of bus i (b_i) represents the number of times a bus i is involved when each pair of buses transmits information along the shortest path between them. The higher the BC, the more frequently the random signal is received, which means that nodes with high b_i are likely to be used for a detour when necessary and have more chances of being vulnerable. The BC protocol is defined as $q_{ji} = b_j/d_{ji}$.

By contrast, the clustering coefficient of node i (c_i) represents the amount of resilience against perturbations. A node with a higher clustering coefficient is less likely to be affected by frequency fluctuations in the electrical grid. Therefore, we define the clustering protocol as $q_{ji} = c_j/d_{ji}$.

Finally, Fiedler's mode is the eigenmode corresponding to the smallest nonzero eigenvalue in the Laplacian matrix and is defined as follows:

$$L_{ij} = \begin{cases} -K_{ij} \cos(\theta_i^* - \theta_j^*) & \text{if } i \neq j, \\ \sum_j K_{ij} \cos(\theta_i^* - \theta_j^*) & \text{if } i = j. \end{cases} \quad (4)$$

The Fiedler mode can be viewed as a set of sensitive buses, with amplitudes representing their sensitivities. Because they respond largely, even to a small external perturbation, it is expected that buses engaging in this mode tend to easily fluctuate around their stable points, which can cause the failure of another bus, eventually leading to a global cascading failure. Therefore, reinforcing the Fiedler mode a_i amplitudes may be an effective method. In general, the mode's amplitude can be negative; therefore, we use it after taking the squared value as $q_{ji} = a_j^2/d_{ji}$.

4.2. GC-PPO protocol

Although classical topological protocols can capture the structural properties of a power grid, they cannot capture the temporal dynamics of a system. Designing a power-dispatch protocol that considers

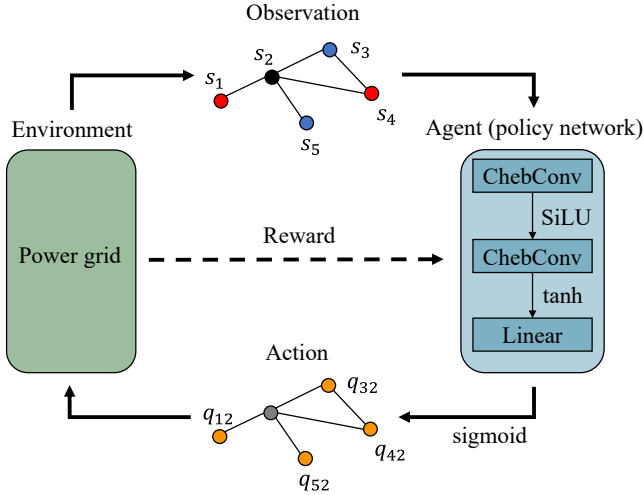


Figure 6: Feedback loop consisting of interactions between the environment (power grid) and agent (policy network). First, the agent observes the steady state and perturbation where generator 2 fails. Second, the agent outputs the action q_{ji} using the observation. This neural network consists of two Chebyshev convolution layers with SiLU and tanh activations and a linear layer with sigmoid activation to constrain the output between 0 and 1. Lastly, the power grid is rebalanced, and the policy network is updated to maximize the reward.

both complex topology and nonlinear dynamics driven by the swing equation (1) is challenging. We propose a reinforcement learning (RL) approach called the Graph Convolutional Proximal Policy Optimization (GC-PPO) protocol to address this issue.

The RL scheme has two main components, the environment and the agent, composed of a power grid and a neural network called a policy network. The agent observes the environment and takes action based on its observations (Fig. 6). The environment provides feedback (reward) on the outcome of the action to the agent to help make better decisions. The agent repeats its interactions with the environment and learns the optimal action for a given observation through trial and error.

Various RL algorithms have been proposed to train agents, and we employ one of the most robust algorithms, the PPO [54]. This was proposed to avoid learning instability caused by fluctuations in the dynamical states by forcing the agent to take actions that are not significantly different from its previous actions during training. In addition, a policy network comprising the agent is implemented using graph convolution to understand the topological properties of the power grid. Notably, after appropriately modifying the structure of a neural network [55, 56], we can

effectively use a single agent for various networks.

GC-PPO training starts with the initial configuration $\{\theta_i, \hat{\theta}_i\}$ in a synchronous state, as defined in Section. 2.3. A perturbation on bus i is given randomly, causing the entire system to lose its power balance. The current states $\{P_i, m_i, \gamma_i, \theta_i, \hat{\theta}_i\}$ of all the nodes, coupling constants $\{K_{ij}\}$, and perturbation δP_i are provided as inputs to the neural network.

For the policy network, we use Chebyshev convolution [55], which considers both the Laplacian of the graph and its higher-order contributions. The first Chebyshev convolution layer with SiLU activation lifts the observations from the environment to a high-dimensional vector (Fig. 6). The second Chebyshev convolution operation, activated by Tanh, computes a high-dimensional node-embedding vector. Next, all the node vectors share the linear and normalization layers and output a scalar value. Finally, the resulting scalar values are passed through a sigmoid function, returning the action $q_{ji} \in [0, 1]$ for the unperturbed bus j . Both Chebyshev layers have identical structures in terms of the number of convolution filters (four filters, considering the power of the Laplacian matrix up to L^3).

However, the output q_{ji} from the policy network is not used directly for power dispatch δP_{ji} as in Eq (3). The GC-PPO protocol behaves differently when the agent is trained and is deployed after training. We consider a Bernoulli distribution with probability q_{ji} for each unperturbed bus j in the training phase. Buses are independently sampled from each distribution to determine whether to participate in the power dispatch. In other words, we set q_{ji} of the nonparticipating bus to zero. This sampling process is crucial for training the agent to perform an optimal action [57, 58]. It enables exploration where the agent experiences diverse situations, even with the same q_{ji} , thereby building robustness into the policy. Note that using the Bernoulli distribution is not the only method to force a probabilistic action. One might sample q_{ji} from a Dirichlet distribution parameterized by $\{\alpha_j\}_{j \in \mathcal{G} \setminus \{i\}}$, where $\alpha_j \in [0, \infty]$. However, the optimization of the agent suffers from the large parameter space of the Dirichlet distribution, making it practically impossible. In contrast, the Bernoulli distribution with a reduced parameter space allows the neural network to be trained efficiently [59].

Meanwhile, deterministic action is preferred in the deployment phase to avoid unpredictable outcomes in critical infrastructures such as power grids. Therefore,

we forgo the sampling process and set a threshold $q_{\text{th}} = 1/10$, where q_{ji} is taken as zero if it is less than q_{th} .

Thus, the output of the policy network is appropriately modified to δP_{ji} following Eq (3). After the power dispatch, the fluctuation Ξ_i of the grid is measured for T seconds. We set $T = 2$ in the training phase to reduce the computational cost, while $T = 10$ is used in the deployment phase to measure the performance of the GC-PPO precisely. As the agent is trained to earn more rewards, $-\Xi_i$ is passed on as a reward, completing one feedback loop. As discussed above, an initial failure can occur on any bus. Therefore, we repeat this feedback loop for every possible i , which we define as an episode.

The agent is trained following the standard PPO method as it progresses through one episode. For each feedback loop initiated by the perturbation of bus i , we compute the probability p_j that the contributing generators will be selected. This is tractable because the generators are sampled in a Bernoulli distribution according to q_{ji} , which is the output of a policy network. Subsequently, the objective of loop L_i is defined as follows:

$$L_i = -\max \left[\frac{p_i}{p'_i} \Xi_i, \text{clip} \left(\frac{p_i}{p'_i}, 1 - \epsilon, 1 + \epsilon \right) \Xi_i \right], \quad (5)$$

where p'_i is the probability calculated using the previous policy network.

The clip function limits the absolute value of the ratio of the probabilities calculated by the current and previous policy networks to $\epsilon = 0.1$, which is given as follows:

$$\text{clip}(x, x_{\min}, x_{\max}) = \begin{cases} x_{\min} & \text{if } x < x_{\min}, \\ x & \text{if } x_{\min} \leq x \leq x_{\max}, \\ x_{\max} & \text{if } x_{\max} < x. \end{cases} \quad (6)$$

Finally, we average the objective L_i over the entire episode and perform a gradient ascent over the averaged objective to update the policy network.

Acknowledgments

This study was supported by the National Research Foundation of Korea No. RS-2023-00279802 (BK) and No. NRF-2022R1C1C1005856 (HK) and a KENTECH Research Grant No. KRG2021-01-007 (BK) and No. KRG2021-01-003 (HK).

The authors declare no conflicts of interest.

The code is available upon request.

References

- [1] O. Smith, O. Cattell, E. Farcot, R. D. O’Dea, K. I. Hopcraft, The effect of renewable energy incorporation on power grid stability and resilience, *Science Advances* 8 (9) (2022) eabj6734.
- [2] P. C. Böttcher, D. Witthaut, L. Rydin Gorjão, Dynamic stability of electric power grids: Tracking the interplay of the network structure, transmission losses, and voltage dynamics, *Chaos: An Interdisciplinary Journal of Nonlinear Science* 32 (5) (2022) 053117. doi:10.1063/5.0082712. URL <https://aip.scitation.org/doi/10.1063/5.0082712>
- [3] D. Florian, C. Michael, B. Francesco, Synchronization in complex oscillator networks and smart grids, *Proceedings of the National Academy of Sciences* 110 (6) (2013) 2005–2010. doi:10.1073/pnas.1212134110. URL <https://doi.org/10.1073/pnas.1212134110>
- [4] A. E. Motter, S. A. Myers, M. Anghel, T. Nishikawa, Spontaneous synchrony in power-grid networks, *Nature Physics* 9 (3) (2013) 191–197. doi:10.1038/nphys2535. URL <https://doi.org/10.1038/nphys2535>
- [5] K. Schmietendorf, J. Peinke, O. Kamps, The impact of turbulent renewable energy production on power grid stability and quality, *The European Physical Journal B* 90 (2017) 1–6.
- [6] A. Ulbig, T. S. Borsche, G. Andersson, Impact of low rotational inertia on power system stability and operation, *IFAC Proceedings Volumes* 47 (3) (2014) 7290–7297.
- [7] O. Edenhofer, R. Pichs-Madruga, Y. Sokona, K. Seyboth, P. Matschoss, S. Kadner, T. Zwickel, P. Eickemeier, G. Hansen, S. Schlomer, C. von Stechow (Eds.), *Renewable Energy Sources and Climate Change Mitigation*, Cambridge University Press, Cambridge, 2011. doi:10.1017/CB09781139151153.
- [8] P. Milan, M. Wächter, J. Peinke, Turbulent Character of Wind Energy, *Physical Review Letters* 110 (13) (2013) 138701. doi:10.1103/PhysRevLett.110.138701.
- [9] M. Anvari, G. Lohmann, M. Wächter, P. Milan, E. Lorenz, D. Heinemann, M. R. R. Tabar, J. Peinke, Short term fluctuations of wind and solar power systems, *New Journal of Physics* 18 (6) (2016) 063027. doi:10.1088/1367-2630/18/6/063027.
- [10] X. Zhang, S. Hallerberg, M. Matthiae, D. Witthaut, M. Timme, Fluctuation-induced distributed resonances in oscillatory networks, *Science Advances* 5 (7) (2019) eaav1027. doi:10.1126/sciadv.aav1027.
- [11] M. Tyloo, P. Jacquod, Primary control effort under fluctuating power generation in realistic high-voltage power networks, *IEEE Control Systems Letters* 5 (3) (2020) 929–934.
- [12] P. J. Menck, J. Heitzig, J. Kurths, H. Joachim Schellnhuber, How dead ends undermine power grid stability, *Nature Communications* 5 (1) (2014) 3969. doi:10.1038/ncomms4969. URL <https://doi.org/10.1038/ncomms4969>
- [13] T. Pesch, H.-J. Allelein, J.-F. Hake, Impacts of the transformation of the german energy system on the transmission grid, *The European Physical Journal Special Topics* 223 (12) (2014) 2561–2575.
- [14] F. Milano, F. Dörfler, G. Hug, D. J. Hill, G. Verbič, Foundations and challenges of low-inertia systems, in: 2018

- power systems computation conference (PSCC), IEEE, 2018, pp. 1–25.
- [15] M. Anvari, F. Hellmann, X. Zhang, Introduction to focus issue: Dynamics of modern power grids, *Chaos: An Interdisciplinary Journal of Nonlinear Science* 30 (6) (2020).
- [16] B. Schäfer, T. Pesch, D. Manik, J. Gollenstede, G. Lin, H.-P. Beck, D. Witthaut, M. Timme, Understanding braess’ paradox in power grids, *Nature Communications* 13 (1) (2022) 5396.
- [17] P. Tielens, D. Van Hertem, The relevance of inertia in power systems, *Renewable and Sustainable Energy Reviews* 55 (2016) 999–1009.
- [18] I. Simonsen, L. Buzna, K. Peters, S. Bornholdt, D. Helbing, Transient dynamics increasing network vulnerability to cascading failures, *Physical Review Letters* 100 (21) (2008) 218701.
- [19] P. J. Menck, J. Heitzig, N. Marwan, J. Kurths, How basin stability complements the linear-stability paradigm, *Nature Physics* 9 (2) (2013) 89–92. doi:10.1038/nphys2516. URL <https://doi.org/10.1038/nphys2516>
- [20] M. Tyloo, R. Delabays, P. Jacquod, Noise-induced desynchronization and stochastic escape from equilibrium in complex networks, *Physical Review E* 99 (6) (2019) 062213.
- [21] J. Hindes, P. Jacquod, I. B. Schwartz, Network desynchronization by non-gaussian fluctuations, *Physical Review E* 100 (5) (2019) 052314.
- [22] B. Schäfer, C. Beck, K. Aihara, D. Witthaut, M. Timme, Non-gaussian power grid frequency fluctuations characterized by lévy-stable laws and superstatistics, *Nature Energy* 3 (2) (2018) 119–126.
- [23] L. Meng, J. Zafar, S. K. Khadem, A. Collinson, K. C. Murchie, F. Coffele, G. M. Burt, Fast frequency response from energy storage systems—a review of grid standards, projects and technical issues, *IEEE Transactions on Smart Grid* 11 (2) (2019) 1566–1581.
- [24] J. Alipoor, Y. Miura, T. Ise, Power system stabilization using virtual synchronous generator with alternating moment of inertia, *IEEE Journal of Emerging and Selected Topics in Power Electronics* 3 (2) (2014) 451–458.
- [25] L. Pagnier, P. Jacquod, Optimal placement of inertia and primary control: A matrix perturbation theory approach, *IEEE Access* 7 (2019) 145889–145900.
- [26] Y. G. Rebours, D. S. Kirschen, M. Trotignon, S. Rossignol, A Survey of Frequency and Voltage Control Ancillary Services—Part I: Technical Features, *IEEE Transactions on Power Systems* 22 (1) (2007) 350–357. doi:10.1109/TPWRS.2006.888963.
- [27] D. Heide, M. Greiner, L. Von Bremen, C. Hoffmann, Reduced storage and balancing needs in a fully renewable european power system with excess wind and solar power generation, *Renewable Energy* 36 (9) (2011) 2515–2523.
- [28] J. Fleer, P. Stenzel, Impact analysis of different operation strategies for battery energy storage systems providing primary control reserve, *Journal of Energy Storage* 8 (2016) 320–338.
- [29] J. Schiffer, R. Ortega, A. Astolfi, J. Raisch, T. Sezi, Conditions for stability of droop-controlled inverter-based microgrids, *Automatica* 50 (10) (2014) 2457–2469.
- [30] H. Taher, S. Olmi, E. Schöll, Enhancing power grid synchronization and stability through time-delayed feedback control, *Physical Review E* 100 (6) (2019) 062306.
- [31] B. Schäfer, D. Witthaut, M. Timme, V. Latora, Dynamically induced cascading failures in power grids, *Nature Communications* 9 (1) (2018) 1975. doi:10.1038/s41467-018-04287-5. URL <https://doi.org/10.1038/s41467-018-04287-5>
- [32] P. C. Böttcher, A. Otto, S. Kettemann, C. Agert, Time delay effects in the control of synchronous electricity grids, *Chaos: An Interdisciplinary Journal of Nonlinear Science* 30 (1) (2020).
- [33] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, D. Hassabis, Mastering the game of Go with deep neural networks and tree search, *Nature* 529 (7587) (2016) 484–489. doi:10.1038/nature16961. URL <https://doi.org/10.1038/nature16961>
- [34] R. S. Sutton, A. G. Barto, Reinforcement learning: An introduction, MIT press, 2018.
- [35] Y.-D. Kwon, J. Choo, B. Kim, I. Yoon, Y. Gwon, S. Min, Pomo: Policy optimization with multiple optima for reinforcement learning, in: H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, H. Lin (Eds.), *Advances in Neural Information Processing Systems*, Vol. 33, Curran Associates, Inc., 2020, pp. 21188–21198.
- [36] A. Mirhoseini, A. Goldie, M. Yazgan, J. W. Jiang, E. Songhori, S. Wang, Y.-J. Lee, E. Johnson, O. Pathak, A. Nazi, J. Pak, A. Tong, K. Srinivasa, W. Hang, E. Tuncer, Q. V. Le, J. Laudon, R. Ho, R. Carpenter, J. Dean, A graph placement methodology for fast chip design, *Nature* 594 (7862) (2021) 207–212. doi:10.1038/s41586-021-03544-w. URL <https://doi.org/10.1038/s41586-021-03544-w>
- [37] J. Degraeve, F. Felici, J. Buchli, M. Neunert, B. Tracey, F. Carpanese, T. Ewalds, R. Hafner, A. Abdolmaleki, D. de las Casas, C. Donner, L. Fritz, C. Galperti, A. Huber, J. Keeling, M. Tsimpoukelli, J. Kay, A. Merle, J.-M. Moret, S. Noury, F. Pesamosca, D. Pfau, O. Sauter, C. Sommariva, S. Coda, B. Duval, A. Fasoli, P. Kohli, K. Kavukcuoglu, D. Hassabis, M. Riedmiller, Magnetic control of tokamak plasmas through deep reinforcement learning, *Nature* 602 (7897) (2022) 414–419. doi:10.1038/s41586-021-04301-9. URL <https://doi.org/10.1038/s41586-021-04301-9>
- [38] J. Alexander, Oscillatory solutions of a model system of nonlinear swing equations, *International Journal of Electrical Power & Energy Systems* 8 (3) (1986) 130–136.
- [39] Q. Qiu, R. Ma, J. Kurths, M. Zhan, Swing equation in power systems: Approximate analytical solution and bifurcation curve estimate, *Chaos: An Interdisciplinary Journal of Nonlinear Science* 30 (1) (2020).
- [40] P. Schultz, J. Heitzig, J. Kurths, A random growth model for power grids and other spatially embedded infrastructure networks, *The European Physical Journal Special Topics* 223 (12) (2014) 2593–2610. doi:10.1140/epjst/e2014-02279-6. URL <https://doi.org/10.1140/epjst/e2014-02279-6>
- [41] F. Dorfler, F. Bullo, Kron reduction of graphs with applications to electrical networks, *IEEE Transactions on Circuits and Systems I: Regular Papers* 60 (1) (2013) 150–163. doi:10.1109/TCSI.2012.2215780.

- [42] D. Witthaut, F. Hellmann, J. Kurths, S. Kettemann, H. Meyer-Ortmanns, M. Timme, Collective nonlinear dynamics and self-organization in decentralized power grids, *Reviews of Modern Physics* 94 (1) (2022) 015005. doi:10.1103/RevModPhys.94.015005.
- [43] C. Mitra, A. Choudhary, S. Sinha, J. Kurths, R. V. Donner, Multiple-node basin stability in complex dynamical networks, *Physical Review E* 95 (3) (2017) 032317. doi:10.1103/PhysRevE.95.032317.
- [44] F. Paganini, E. Mallada, Global performance metrics for synchronization of heterogeneously rated power systems: The role of machine models and inertia, in: 2017 55th Annual Allerton Conference on Communication, Control, and Computing (Allerton), IEEE, 2017, pp. 324–331. doi:10.1109/ALLERTON.2017.8262755.
- [45] C. Nauck, M. Lindner, K. Schürholt, H. Zhang, P. Schultz, J. Kurths, I. Isenhardt, F. Hellmann, Predicting basin stability of power grids using graph neural networks, *New Journal of Physics* 24 (4) (2022) 043041. doi:10.1088/1367-2630/ac54c9.
- [46] S.-G. Yang, B. J. Kim, S.-W. Son, H. Kim, Power-grid stability predictions using transferable machine learning, *Chaos: An Interdisciplinary Journal of Nonlinear Science* 31 (12) (2021) 123127. doi:10.1063/5.0058001.
- [47] C. Nauck, M. Lindner, K. Schürholt, F. Hellmann, Toward dynamic stability assessment of power grid topologies using graph neural networks, *Chaos: An Interdisciplinary Journal of Nonlinear Science* 33 (10) (2023).
- [48] B. Jhun, H. Choi, Y. Lee, J. Lee, C. H. Kim, B. Kahng, Prediction and mitigation of nonlocal cascading failures using graph neural networks, *Chaos: An Interdisciplinary Journal of Nonlinear Science* 33 (1) (2023) 013115. doi:10.1063/5.0107420.
- [49] L. C. Freeman, A Set of Measures of Centrality Based on Betweenness, *Sociometry* 40 (1) (1977) 35–41. doi:10.2307/3033543. URL <http://www.jstor.org/stable/3033543>
- [50] L. C. Freeman, Centrality in social networks conceptual clarification, *Social Networks* 1 (3) (1978) 215–239. doi:https://doi.org/10.1016/0378-8733(78)90021-7.
- [51] D. J. Watts, S. H. Strogatz, Collective dynamics of ‘small-world’ networks, *Nature* 393 (6684) (1998) 440–442.
- [52] L. Pagnier, P. Jacquod, Inertia location and slow network modes determine disturbance propagation in large-scale power grids, *PLOS ONE* 14 (3) (2019) e0213550. doi:10.1371/journal.pone.0213550.
- [53] Z. Wang, A. Scaglione, R. J. Thomas, Electrical centrality measures for electric power grid vulnerability analysis, in: 49th IEEE Conference on Decision and Control (CDC), 2010, pp. 5792–5797. doi:10.1109/CDC.2010.5717964.
- [54] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, O. Klimov, Proximal policy optimization algorithms (2017). arXiv:1707.06347.
- [55] M. Defferrard, X. Bresson, P. Vandergheynst, Convolutional neural networks on graphs with fast localized spectral filtering, *Advances in Neural Information Processing Systems* 29 (2016).
- [56] J. Jiang, C. Dun, T. Huang, Z. Lu, Graph convolutional reinforcement learning, in: International Conference on Learning Representations, 2020. URL <https://openreview.net/forum?id=HkxdQkSYDB>
- [57] S. Ishii, W. Yoshida, J. Yoshimoto, Control of exploitation–exploration meta-parameter in reinforcement learning, *Neural Networks* 15 (4-6) (2002) 665–687.
- [58] M. Castronovo, F. Maes, R. Fonteneau, D. Ernst, Learning exploration/exploitation strategies for single trajectory reinforcement learning, in: European Workshop on Reinforcement Learning, PMLR, 2013, pp. 1–10.
- [59] B. Li, Z. Wei, J. Wu, S. Yu, T. Zhang, C. Zhu, D. Zheng, W. Guo, C. Zhao, J. Zhang, Machine learning-enabled globally guaranteed evolutionary computation, *Nature Machine Intelligence* 5 (4) (2023) 457–467. doi:10.1038/s42256-023-00642-4. URL <https://doi.org/10.1038/s42256-023-00642-4>