

---

# Efficient Management of Day-Ahead Energy Markets via Multi-Agent Reinforcement Learning - a Hybrid Model Case Study

---

Matan Levy<sup>1</sup>, Itay Segev<sup>1</sup>, Alexander Tuisov<sup>1</sup>, Sarah Keren<sup>1</sup>

<sup>1</sup> Technion-Israel Institute of Technology

{matanlevy,itaysegev,alexandt}@campus.technion.ac.il, sarahk@technion.ac.il

## Abstract

This study examines the optimization of day-ahead hybrid electricity markets. The shift from centralized systems to public-private models introduces many challenges, including the introduction of independent market players and *renewable energy sources* (RESs). A formal model of market participants' behavior is developed, and a *multi-agent reinforcement learning* (MARL) framework is proposed to optimize system operator strategies, incorporating dynamic pricing and dispatch scheduling to reduce operational costs, ensure stability, and align market incentives. A new and adaptable simulation environment, compatible with state-of-the-art methods, is presented. Evaluations in increasingly complex settings demonstrate the efficacy of our framework in managing the complexities of modern electricity markets.

## 1 Introduction

This work addresses the day-ahead optimization of an electricity market<sup>1</sup> undergoing significant structural transformation. Historically centralized and government-controlled, the increasing integration of *renewable energy sources* (RESs) and the advancements in data collection technologies are transitioning the market into a complex public-private hybrid model. This presents substantial challenges and the need to deal with a highly uncertain operational and regulatory environment [1].

To demonstrate some of the challenges involved in managing current energy systems, consider a day-ahead market in which the *independent system operator* (ISO) aims to optimize electricity generation based on forecasted demand, generation costs, and grid constraints. The resulting decisions, made 24 hours in advance, specify the amount of electricity to be produced, the prices, and the allocation of reserve capacity, i.e., the ability to generate additional power at short notice, often at high environmental costs, in the event of generation failures or unexpected demand surges.

Adapting the day-ahead market to today's energy systems requires accounting for the variability and limited controllability of increasingly heterogeneous *grid-edge agents*, denoted hereon as **GEAgents**, particularly those with local generation and storage capabilities. For example, a household with a photovoltaic (PV) unit and a battery can autonomously optimize its energy storage policy, learning when to store energy, when to consume it, and when to trade with the grid to maximize economic benefits. While such behavior may improve individual utility, it introduces significant uncertainty into aggregate demand forecasts and can destabilize the system, especially under sudden shifts in consumption or generation patterns. At the same time, these distributed resources can enhance efficiency and resilience by shaving peaks, supplying energy, and reducing the amount of centrally dispatched generation required.

---

<sup>1</sup>For anonymity reasons, the specific market under consideration is not disclosed.

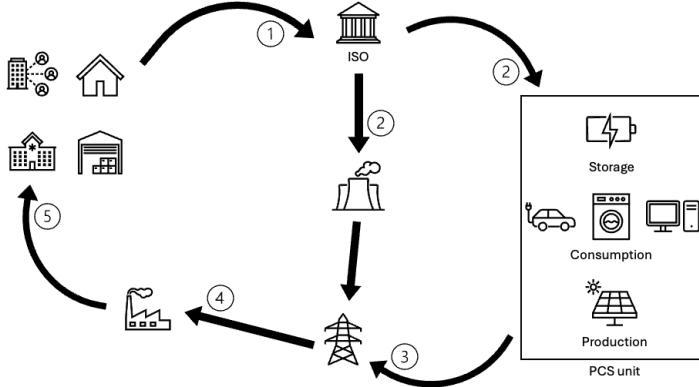


Figure 1: The day-ahead control cycle that repeats every 30-minutes: (1) ISO receives realized demand for the current time step. (2) ISO posts real-time buy/sell tariffs and issues dispatch directives to the controlled generators (3) GEAgents buy/sell power (4) If needed, peaker reserves are dispatched or curtailment is performed (5) Balanced power flows to consumers.

To address these challenges, the ISO adjusts electricity production plan, or *dispatch*, and feed-in and sell prices to influence independent market participants and align their behavior with grid operational objectives. Additionally, it retains access to reserves and peaking power plants, which can be activated to address unmet demand, ensuring both system stability and operational efficiency. The problem the ISO faces is thus one of cost optimization while satisfying the demand in the presence of strategic market players that aim to maximize their own profits. The scale and complexity of the problem make data-driven approaches, such as *reinforcement learning* (RL), especially suitable [2].

We make three key contributions. First, we build a ***multi-agent reinforcement learning (MARL)*** model that captures the incentives and rational decision-making of independent market participants. Leveraging these models, we then study the ISO’s optimization problem under various assumptions, revealing how each setting shapes optimal dispatch and pricing policies. Finally, we offer a configurable, open-source grid simulator that supports diverse topologies and uncertainty patterns. Experiments across increasingly complex settings demonstrate that RL-driven agents can jointly optimise participant and ISO strategies, highlighting the promise of MARL for modern energy-market design.

## 2 Background and Related Work

Reinforcement Learning (RL) is a learning paradigm where an agent learns optimal behavior by interacting with an environment and receiving rewards or penalties for its actions [3]. Multi-agent reinforcement learning (MARL) extends RL to scenarios involving multiple autonomous agents that concurrently learn and make decisions within a shared or partially shared environment [4]. Each agent aims to maximize its own utility (typically measured as accumulated reward), but its actions can influence both its own outcomes and the outcomes of other agents, leading to complex emergent behaviors and the need for coordination and cooperation (see Appendix A for more detail).

The most common MARL model is the *stochastic game* (SG) (also known as emMarkov game or multi-agent MDP) [5] defined as a tuple  $\langle \mathcal{S}, \mathcal{A} = \{\mathcal{A}_i\}_{i=1}^n, \mathcal{T}, \mathcal{R} = \{\mathcal{R}_i\}_{i=1}^n, \gamma \rangle$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the *joint action space* with  $\mathcal{A}_i$  as the  $i^{th}$  agent action space s.t.  $a \triangleq (a_1, a_2, \dots, a_n)$  for  $a \in \mathcal{A}$ ,  $\mathcal{T} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  is the transition probability function  $\mathcal{T}(s', a, s)$  such that  $\forall s \in \mathcal{S}, \forall a \in \mathcal{A} : \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') = 1$ ,  $\mathcal{R}$  is the *joint reward function* with  $\mathcal{R}_i : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  as the  $i^{th}$  agent reward function, and  $\gamma \in [0, 1]$  is the discount factor. A solution is a joint policy  $\pi \triangleq (\pi_1, \dots, \pi_n)$  associating each agent with policy  $\pi_i : \mathcal{S} \times \mathcal{A}_i \rightarrow [0, 1]$  that specifies the probability of agent  $i$  taking an action at a given state. The joint policy should achieve certain conditions on the expected returns yielded to agents (e.g., Nash equilibrium) [4]. The value (utility) function  $V_i^\pi(s)$  denotes the expected cumulative discounted reward agent  $i$  receives when starting in state  $s$  and the agents follow joint policy  $\pi$  thereafter. The action-value function or Q-value  $Q_i^\pi(s, a)$  extends this

notion by quantifying the expected value when performing  $a$  in  $s$ , and then continuing according to  $\pi$ . This general definition captures a variety of interactions and relationships that can exist between agents in collaborative, competitive, and mixed-incentive MARL settings.

MARL is particularly suitable for modeling energy systems and networks, since they are inherently multi-agent environments composed of diverse, distributed, and strategically autonomous entities, such as grid-edge components, utility companies, system operators, and market participants [1]. These entities have different objectives, interact over shared physical and economic infrastructures, and must respond dynamically to system conditions, prices, and regulations. MARL provides a natural framework to model these interactions, enabling agents to learn adaptive policies, coordinate under uncertainty, and reason about both cooperative and competitive dynamics. Moreover, its ability to simulate emergent behavior and explore decentralized strategies makes it a powerful tool for both designing and analyzing modern energy systems.

Applications of RL and MARL in energy markets often assume a single, all-knowing controller optimizing the entire system. In such formulations, a central agent (analogous to an ISO) directly controls all generation and storage decisions using global information and perfect foresight, an assumption that is unattainable in practice. These centralized optimization models can yield system-level insights but cannot capture the strategic, profit-driven behavior of individual market participants [6, 2]. Moreover, as modern grids grow more heterogeneous and stochastic with high renewable penetration, a monolithic control scheme becomes impractical[7]. Recent studies emphasize that managing numerous distributed resources under uncertainty requires moving beyond one-size-fits-all control toward more decentralized decision-making structures [8, 9].

On the other end of the spectrum, many RL-based models use a fully decentralized approach in which each market participant (e.g. a storage unit owner or consumer) acts independently. In these formulations, multiple RL agents learn their own policies (for bidding, charging, discharging, etc.) based on price signals or local observations, without a central coordinator explicitly optimizing the whole system[10]. This bottom-up approach reflects competitive markets by giving each market player its own profit-maximizing RL agent[11, 12, 13]. However, purely decentralized models typically assume the market rules or prices are exogenous or fixed [1, 14, 2]. In our model, the ISO acts as an active participant and directly shapes the market dynamics. Related efforts on dynamic dispatch and end-to-end RL in energy systems include [15, 16], and comprehensive overviews of RL for power systems can be found in [14].

From an algorithmic standpoint, our challenge is to address the strategic and tightly coupled interplay between a central ISO and price-responsive market participants. Most prior approaches either employ fully centralized optimization, which ignores competitive behavior, or simulate independent agents interacting with static ISO actions. Consequently, hybrid markets, where a dynamic, learning-enabled ISO coexists with autonomous market agents, remain underexplored. In our framework, the ISO continually adjusts dispatch and pricing signals, while market agents react strategically to maximize their own profit. Capturing this two-way interaction is essential for realistic market modeling, yet most RL studies to date have only touched on limited aspects of this ISO-agent feedback loop [6, 17]. The scarcity of work in this hybrid paradigm highlights the potential of methods like the one we proposed for integrating a central coordinator’s adaptive decisions with the learning-based responses of individual market players and for optimizing modern energy systems.

### 3 Energy Market Dynamics

Historically, the energy market comprised three principal components: power producers (e.g., power plants), power consumers (industrial and residential), and the ISO, responsible for market management and coordination. The producers typically used conventional coal-based generation and were either units under the full control of the ISO, or independent units that participated in the market but were regulated and bound by production agreements made for different temporal horizons.

In a typical *day-ahead market*, as depicted in Figure 1, the ISO predicts the following day’s power demand (electricity consumption) and issues a *dispatch*, a production schedule, while considering operational constraints and generation costs. In addition to the generation of the predicted, or *nominal* demand, the ISO also manages the *reserve*, which sets a backup production capability for each time step. In real-time, the ISO is tasked with continuously maintaining a balance between demand and supply. If there is a surplus, energy is discharged, or *curtailed*. If production determined by the

dispatch is not enough to cover the *realized demand*, reserves, which are more flexible but also more expensive and polluting, are deployed. Producers are then compensated based on the System Marginal Price (SMP) mechanism, calculated as the marginal cost of producing the final unit of energy required to satisfy system demand, based on the least-cost dispatch solution. In this work, we abstract the dispatch details and consider only the total amount and cost of power produced at each timestamp (see Appendix B and C for details on market dynamics and SMP computation, respectively).

Recent power-market reforms have introduced independent grid-edge agents (GEAgents), from private utilities to smart homes, alongside traditional producers and consumers. In these markets, each GEAgent operates a **Production-Consumption-Storage unit (PCS-unit)**, which may produce (e.g., via PV), consume (e.g., via electrical appliances), and store (e.g., via a battery) energy. Unlike traditional controlled producers, they are not legally required to adhere to dispatch instructions and may buy from or sell energy to the grid at will to maximize their profits, overlooking the resulting instability. Therefore, modeling the GEAgents' behaviors and strategies is essential for the ISO's planning. Since we assume GEAgents are rational, the natural way for the ISO to align player incentives with stability constraints and efficiency objectives is via pricing. Thus, with the aim of minimizing total costs for the ISO (thus the taxpayers) while satisfying the supply and demand balance, the dispatch, denoted  $\Delta_t$ , and selling and feed-in prices, denoted  $\xi_t$  and  $\phi_t$ , respectively, for each time  $t$ , are the primary tools for market control. In what follows, we analyze the ISO's optimization problem under increasingly complex market regimes, highlighting how dispatch and pricing jointly control system stability and efficiency.

In the deterministic setting, fully formulated in Appendix B, the ISO receives at the beginning of each episode the nominal production and reserve capabilities and costs for market participants, as well as the demand for all time steps in the horizon  $T$ . Based on this information and the operational constraints, it determines the scheduled  $\Delta_t$  and prices  $\xi_t(\cdot)$ ,  $\phi_t(\cdot)$  for all timestamps  $t \in [T]$  to minimize total costs. Formally,

$$\min C^{\text{total}} = \min \left[ C^{\text{dispatch}} + \sum_{t=1}^T C_t^{\text{online}} \right] \quad (\text{Deterministic ISO Objective})$$

where  $C^{\text{dispatch}}$  is the total dispatch cost for the complete episode, and  $C_t^{\text{online}}$  is the online cost (including reserve cost) for time  $t$ .

Since all information is given in advance, the GEAgent can also compute its policy at the beginning of each episode and decide how much power to buy from ( $P_t^b$ ), and sell to ( $P_t^s$ ) the grid at every timestamp  $t$  to maximize its total revenue under its operational constraints. Formally:

$$\max \sum_{t=1}^T (\phi_t P_t^s - \xi_t P_t^b) \quad (\text{Deterministic GEAgent Objective})$$

In a stochastic extension of this setting, we account for the inability to exactly predict demand and production. In this case, it may be possible to estimate these distributions from historical data and observations using machine learning methods to improve decision-making under these forms of uncertainty. In this setting, fully formulated in Appendix B, the min and max objectives of the ISO and GEAgents are replaced by an expectation-based optimization.

**Accounting for Strategic Demand:** In modern energy systems, demand is not only stochastic but also strategic since GEAgents can intelligently manage the operation of devices and energy resources, in response to system-level signals. This *demand (load) flexibility* is reshaping energy markets by introducing new ways to contribute to their efficient and stable operation [18, 1]. However, this shift also introduces challenges such as increased system complexity, uncertainty in demand forecasting, and the need for regulatory mechanisms to ensure fair and reliable participation.

In this extended setting, the ISO needs to determine the selling price  $\xi_t$  and feed-in prices  $\phi_t$  for each  $t$  according to the demand  $D_t$  at time  $t$  while accounting for the GEAgents' ability to sell, buy, and store power. From the perspective of the GEAgent, the price signals  $\xi_t(P_t^s, P_t^b, \dots)$  are exogenous signals set by the ISO, but they depend on the GEAgents' sales  $P_t^s$  and purchases  $P_t^b$  and other variables. This coupling results in a feedback mechanism where the player's actions influence the prices, and the prices, in turn affect the player's actions. This introduces a game-theoretic dimension where the GEAgents' decisions are influenced by the ISO's pricing strategy and vice versa.

Formally, the GEAgent's input includes all the parameters that were relevant for the deterministic and stochastic settings, including the expected demand  $l_t$  and production  $g_t$  at time  $t$ . A key difference is that the selling price  $\xi_t$  and feed-in prices  $\phi_t$  can be set either in advance or, depending on regulation, dynamically, in response to the market state. The objective of the GEAgent is now:

$$\max_{P_t^b, P_t^s} \mathbb{E}_{l_t, g_t} \left[ \sum_{t=1}^T (\phi_t(P_t^s, P_t^b, \dots) - \xi_t(P_t^s, P_t^b, \dots)) \right] \quad (\text{Strategic Player Objective})$$

From the perspective of the ISO, as in the stochastic settings, it receives at the beginning of each episode (day) all the information about the GEAgents and the controlled producers and needs to determine the scheduled amount of production  $\Delta_t$  for each timestamp. However, it is crucial to distinguish between two components of the demand. The **nominal demand** refers to the exogenous, inelastic portion of load that remains unaffected by local control strategies, real-time market incentives, or variations in renewable generation. In contrast **flexible demand**, refers to the portion of demand that can be adjusted in time, quantity, or pattern in response to external signals, such as price changes, grid conditions, or availability of renewable energy.

Since the ISO cannot loyally model the demand without considering the strategic nature of the GEAgents, optimization methods that are appropriate for deterministic and stochastic settings won't work here. Thus, as we specify in the next section, we model the market participants as RL agents.

## 4 The Energy Market as MARL

In modeling modern power systems using MARL, it is essential to account for multiple interacting perspectives. These include the physical constraints of the grid (e.g., stability limits), agent-level decision processes under partial observability, and the heterogeneity of demand profiles encompassing both nominal and flexible demand. Effective models must also incorporate market and pricing signals that influence agent behavior, and the temporal-spatial scalability required for real-world deployment. While these considerations are crucial for realistically and robustly capturing decentralized control strategies in complex energy environments, they also pose significant challenges to preserving the underlying Markovian structure that traditional agent-based decision models rely on.

### 4.1 Formal Model

Through the lens of RL, the ISO aims to learn an optimal policy that balances overall system efficiency with the mitigation of risk, such as insufficient power supply and grid instability. Simultaneously, GEAgents seek to maximize their individual utility in response to market signals, subject to their own operational constraints and preferences. We formally model this decentralized setting as a Markov game (see Section 2), involving two types of agents: the ISO, and the GEAgents. A key characteristic of the setting we aim to model is that the states, actions, and rewards are relatively straightforward to define. The complexity of solving this setting arises in modeling the joint transition function: the next system state and its stability depend on the actions performed by all agents.

#### Modeling the ISO

- **State Space  $\mathcal{S}$ :** Every time step  $t$ , typically representing a half-hour interval, the system state is associated with a vector  $s_t \in \mathcal{S}$  that specifies operational factors that may affect decision-making. For the ISO this includes the system-level demand forecast  $\hat{D}_t$  for the specified horizon, the system-level realized demand  $D_t$  for the current time step, supply capacities, storage states, etc. It may also include factors that indicate the stability state of the system, for example, whether the supply-demand balance is violated.
- **Action Space  $\mathcal{A}$ :** The ISO's actions include the dispatch directives  $\Delta_t$  and setting the sell prices  $\xi_t(\cdot)$  and buy prices  $\phi_t(\cdot)$  for each time step. In real-time, the ISO also activates reserves and curtails power if needed, but since we assume these actions are dictated by the state, they are not modeled as actions.

We support two types of pricing mechanisms. In a day-ahead pricing regime, the ISO makes the prices public at  $t = 0$  while in an online pricing setting, the ISO can dynamically set prices in response to market signals. We discuss several pricing mechanisms and their characteristics, including the benefits of applying quadratic pricing, in Section 5.

- **Reward Function  $\mathcal{R}$ :** We introduce two mismatch weights  $\eta_o, \eta_u > 0$  with  $\eta_u > \eta_o$ , so that at each time  $t$  the ISO’s reward is

$$\mathcal{R}_t(\eta_o, \eta_u) = -(C_t^{\text{dispatch}} + C_t^{\text{online}}) - \eta_o [\max\{0, \Delta_t - D_t\}]^2 - \eta_u [\max\{0, D_t - \Delta_t\}]^2.$$

Under-dispatch ( $\Delta_t < D_t$ ) therefore incurs the larger penalty  $\eta_u$ , while over-dispatch uses  $\eta_o$ .

## Modeling the GEAgnets

- **State Space  $\mathcal{S}$ :** Each GEAgnet is associated with a PCS-unit for which the state includes its local information (e.g., state-of-charge) as well as the price signals advertised by the ISO.
- **Action Space  $\mathcal{A}$ :** Modern GEAgnets have significant decision-making autonomy, allowing them to choose how much energy to store, consume, or sell based on their local goals, capabilities, and constraints. We assume the GEAgnet sees the current prices and local state at the start of each iteration before deciding how to act. Also, both generation and consumption are non-controllable, corresponding to PV-based generation and consumption from appliances. This means that generation and production are exogenous to the agent and are governed by stochastic processes, and the only decision variables are the charge and discharge actions, which may have stochastic effects.
- **Reward Function  $\mathcal{R}$ :** For each GEAgnet  $i$ , the step-wise reward is the net revenue obtained by trading with the grid:  $\mathcal{R}_t^i = \phi_t(P_t^s, P_t^b) - \xi_t(P_t^s, P_t^b)$ . Maximising the cumulative sum of  $\mathcal{R}_t^i$  over the horizon is equivalent to the strategic objective stated in Section 4.

**Joint Transition Function  $\mathcal{T}$ :** Unlike a single-agent MDP, the Markov game framework allows each agent’s choice (including how GEAgnets respond to prices or storage opportunities) to influence the next state. As mentioned above, the difficulty of modeling the transition function is at the core of the challenge. In general, the transition function can be decoupled into two parts. The physical dynamics capture the dynamics of the electrical network. For example, when a charge or discharge action is performed, the battery dynamics must obey its physical constraints. In contrast, market dynamics capture the interactions between agents. These strategic decisions create a coupled system where each agent’s payoff depends on the actions of others. In principle, the Markov transition function must capture all aspects of the dynamics, but writing a closed-form for these different layers is hopeless. Instead, we create the Energy-Net simulator (Section 6) to maintain the physics and book-keeping, and *learn* directly from roll-outs. This side-steps the need for explicit modeling of the complex dynamics and allows extracting value functions and policies using deep neural networks, rather than from first principles.

**Episode:** As is typical in the day-ahead market, at the beginning of each episode (timestep  $t = 0$ ) the ISO receives the predicted demand  $\hat{D}_t$  for the next 48 half-hour intervals. It also receives the production and reserve capacities of its controlled units, the prices of each generated unit, and other information that might be relevant (i.e., weather forecast, special events, etc.). If day-ahead pricing is applied, the ISO sets and advertises the  $\xi_t(\cdot)$  and feed-in tariff  $\phi_t(\cdot)$  for the whole episode. Otherwise, online pricing is applied. This iterative process continues until the end of the planning horizon. The full cycle is described in Appendix G and depicted in Figure 1.

## 5 Solution Approaches

The MARL formulation described in Section 4 provides an abstraction that captures the strategic, price-driven interactions that typify modern hybrid power systems. In this section, we present solution approaches that can be adopted by the market participants. Importantly, while our main challenge is in computing optimal market management approaches for the ISO, we must equip the GEAgnets with the strongest policies to guarantee the ISO can predict their response to different price signals.

In principle, the deterministic and stochastic formulations described in Section 3 can be solved using state-space and dynamic programming methods, respectively (see Appendix D for an example formulation). Even if distributions are not fully known, it may be possible to learn them from data. Nevertheless, such methods are not appropriate to our problem, which is inherently challenging due

to the agents’ ability to strategically adapt their behavior and due to the dual-action learning structure, which operates across different time frames.

A specific challenge is that pricing may be dynamic and set at every time step, while the  $\Delta_t$  action for each time step  $t$  is decided at the beginning of each episode. This temporal disparity adds a layer of complexity, as the reward for a  $\Delta$  action is reflected only at the end of the episode. Moreover, determining  $\Delta$  is a demanding task because it involves generating a time series output that must account for dynamic market conditions, which are influenced by behaviors of market participants. A further complication arises from the interdependence of these actions. Dispatch decisions are influenced by the market agents’ responses to price signals, while optimal pricing strategies depend on real-time  $D_t$  and  $\Delta_t$  outcomes.

Because the game is sequential (ISO first, GEAgent second) and highly non-linear, we iteratively train each of the policies with deep RL for continuous control in an online regime. If the agents’ policies converge, it is toward a *practical* equilibrium in function-approximation space rather than a formal Nash point. In Section 7, we empirically examine this using our simulated environment described in the next section.

There are several abstractions that we can use to facilitate computation. One option is to make the problem easier is by abstracting away the dispatch optimization, which we denote as **dispatch abstraction**. In this simplified model the ISO has only control over the prices, and we assume that the ISO production  $\Delta_t$  is fixed to be equal to the predicted demand  $\hat{D}_t$ .

**Quadratic Pricing:** We employ two pricing regimes, online dynamic and day-ahead tariffs. In settings restricted to day-ahead pricing, *quadratic pricing* allows the ISO to influence consumption and injection patterns through price curvature. Following [19], we impose a superlinear surcharge on purchases and a sublinear bonus on feed-in:

$$\xi_t = \alpha_0 + \alpha_1 P_t^b + \alpha_2 [P_t^b]^2, \quad \phi_t = \beta_0 + \beta_1 P_t^s + \beta_2 \sqrt{P_t^s},$$

where the six coefficients  $(\alpha_0, \alpha_1, \alpha_2, \beta_0, \beta_1, \beta_2)$  are fixed at the episode’s outset for the subsequent  $T$  time steps. The superlinear term steepens the marginal purchase price, thereby discouraging demand spikes and reducing reliance on peaker reserves, while the sublinear feed-in adjustment tempers incentives for excessive injections, promoting smoother system operation (see Appendix E for full details and detailed examples).

## 6 The Energy-Net Simulator

In spite of a variety of simulators that currently exist [20, 21, 22, 23], there is no current framework that allows modeling the complex structure we want to account for and that is designed to work with off-the-shelf RL and MARL methods. We therefore develop a novel simulator, Energy-Net, that we will use to examine our proposed solutions. Energy-Net is a modular, discrete-time simulator of a hybrid electricity market. The environment we develop is flexible and adaptable, and can be used to accommodate different system configurations. At the core of the design of the software is a decoupling between the physical dynamics of the electrical system and the strategic agents, i.e., it is built around a strict *physics–agent split*. A high-fidelity physical core advances loads, renewables, batteries, and reserves, while the ISO and GEAgents interact only through a Gym-style `step()` interface. This design (i) lets us plug in any off-the-shelf RL algorithms without touching the power-system code, (ii) isolates market rules in a single controller module, and (iii) ensures that learned policies can affect the grid *only* via explicit levers, prices and dispatch tweaks, thus preserving physical realism while streamlining experimentation.

Building on the formal setting introduced in Section 4, Energy-Net instantiates the 24-hour day-ahead electricity market. A single simulation episode therefore comprises  $T$  uniform intervals of length  $\Delta t$  (in our experiments  $T=48$  and  $\Delta t=30\text{min}$ ), together covering one 24-hour operational horizon. At each step  $t \in \{1, \dots, T\}$  the environment reveals the current forecast and grid state to the agents, applies their actions, propagates the physical dynamics, and returns next-state observations and rewards through the standard Gym `step` interface. See Appendix H for the full details.

## 7 Empirical Evaluation

The objective of our empirical evaluation is to assess the benefit of using our MARL formulation to optimize the policy of the ISO. For this, we use our Energy-Net environment to model and simulate the day-ahead electricity market<sup>2</sup>.

**Setup** We evaluate our formulation from Section 4 and pricing schemes from Section 5 under a variety of scenarios. As discussed in Section 6, Energy-Net cleanly separates physical dynamics from agent logic. This allows us to stage the empirical study in three escalating phases of coordination for the ISO and GEAgents. First, in ISO-Dispatch, we trained and evaluated the ISO in isolation; all GEAgents were disabled, so the operator optimised its dispatch  $\Delta_t$  under a stochastic yet *non-strategic* demand profile. Next, we enabled a PCS-unit<sup>3</sup> with a fixed, pre-defined charging trajectory and retrained the ISO, thereby quantifying the benefit of price coordination when storage is present but *non-adaptive*. We examined this setting with two pricing mechanisms: *online linear*, denoted ISO-L, and *quadratic*, denoted ISO-Q. We then allowed *both* agents to learn concurrently: the ISO tunes its real-time dispatch and tariffs, while the PCS-unit adapts its behavior to these market signals. In settings Joint-Storage-L and Joint-Storage-Q we examined the online and linear pricing, respectively, for a storage-only GEAgent, while in Joint-PCS-L and Joint-PCS-Q, we added production and consumption capabilities (see Appendix I for the full details of the setup). For each episode, we sample the *realized* demand from a Gaussian noise induced predicted demand for each time step  $t$ , and, when relevant, the realized load and production for the PCS-units. (see Table 4 in the appendix for a full description). We ran each training phase for 40 iterations with 4800 time steps each (1000 days) and was evaluated for 20 times. All settings were run using the same demand pattern and performance parameters described in Appendix I with Allocated resources of : 10 cores of Intel(R) Xeon(R) CPU E5-2683 v4 @ 2.10 GHz and 1 × NVIDIA GeForce GPU (12 GB). .

**Results** Due to space constraints, we present our full results in Appendix J and show here only our key findings. Our focus is on optimizing the ISO and measuring its ability to avoid failure and minimize cost, thus preferring to exploit renewable energy generated by the GEAgents and avoiding usage of reserves as much as possible. We therefore present in Table 6 the average energy usage achieved for all multi-agent settings compared to baseline ISO-Dispatch. To fully appreciate the effect of each agent setup, we present a breakdown of the total energy in MWh into three components: dispatch, reserve, and exchange (variance values in parentheses).

Results show that for settings ISO-L and ISO-Q, in which the GEAgent is fixed, the ISO manages to learn to exploit the power generated by the GEAgents instead of the reserves. In contrast, in Joint-Storage-L and Joint-Storage-Q, with a storage-only GEAgent the PCS-unit energy does not contribute to the overall efficiency. Instead, it increases the amount the ISO produces via dispatch to maintain stability. In Appendix J we show how this effect can be mitigated with different cost coefficients. Finally, for the complete setup of Joint-PCS-L and Joint-PCS-Q, where the GEAgents have consumption and production capabilities, we see a minimization of the reserve with quadratic pricing. To further demonstrate GEAgents contribution, Figure 7 depicts an episode from the Joint-PCS-L and Joint-PCS-Q settings. The difference between the dashed black line and the blue line (realized demand) represents the gap between the nominal predicted demand and the realized demand. The dispatch is represented by the light blue bars, while the total demand, including the flexible load of the GEAgents is depicted by the red line (total demand). As demonstrated in the figure, the reserve activation happens when the red line is *above* the dispatch bars, which is to be avoided. Overall, our experiments show that while fixed-generation players (ISO-L and ISO-Q) enable the ISO to substitute market output for reserves and storage-only players (Joint-Storage-L and Joint-Storage-Q) can unintentionally boost dispatch, it is only the combined consumption-production scenario (Joint-PCS-L) under a quadratic day-ahead tariff that suppresses reserve activation and maximizes system efficiency.

---

<sup>2</sup>To respect the blind review process, our code base and complete results are in the supplementary material. All will be made public after acceptance.

<sup>3</sup>Additional units can be added using the same interface; for clarity, we use one aggregated unit.

Table 1: Episode–total *energy* in MWh breakdown across scenarios.

Scenario	Dispatch	Reserve	Exchange
ISO-Dispatch	$7229.86 \pm 38.29$	$249.41 \pm 5.04$	NA
ISO-L	$7282.34 \pm 50.89$	$176.05 \pm 19.07$	$800 \pm 0$
ISO-Q	$7506.98 \pm 35.02$	$121.07 \pm 3.78$	$800 \pm 0$
Joint-Storage-L	$8126.13 \pm 1.07$	$148 \pm 0.94$	$0 \pm 0$
Joint-Storage-Q	$8126.21 \pm 1.01$	$148 \pm 1.06$	$0 \pm 0$
Joint-PCS-L	$7322.44 \pm 36.02$	$168.47 \pm 4.14$	$442.14 \pm 9.61$
Joint-PCS-Q	$7450.62 \pm 36.43$	$117 \pm 2.04$	$324 \pm 8.40$

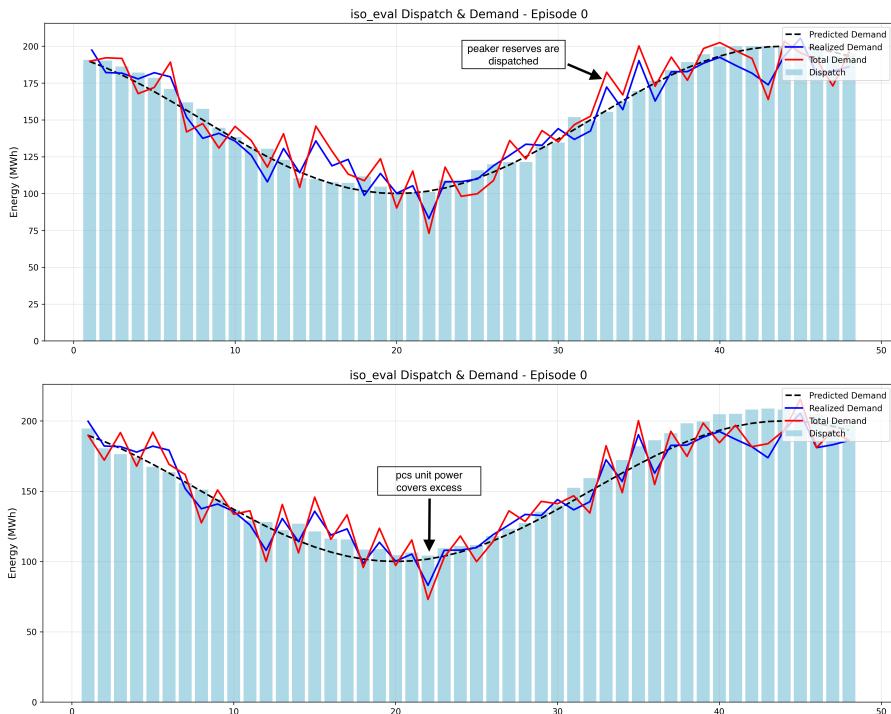


Figure 2: Episode-level dispatch and realized demand under scenario Joint-PCS-L (online linear pricing) at the top, and scenario Joint-PCS-Q (quadratic pricing) at the bottom.

## 8 Conclusion

We demonstrate the benefit of modeling modern power systems MARL in which physical grid constraints, market signals, and heterogeneous agent behaviors interact in tightly coupled feedback loops. We design our framework to capture both nominal and flexible demand, and enable realistic and robust evaluation of decentralized control strategies and pricing mechanisms using a new simulation environment we developed. Our results show that strategically coordinated ISO policies working with price-responsive grid-edge agents can reduce reserve requirements and carbon intensity.

Together with these achievements, our experiments reveal the fragility of current deep-RL policies: modest forecasting errors can lead to supply shortfalls or excessive generation. Addressing this brittleness remains a key research priority. Another challenge lies in scaling the approach operational grids. This will require hierarchical or federated MARL architectures and hardware-in-the-loop testing. Finally, while algorithmic coordination can reduce reserve usage and lower tariffs, distribution benefits are unlikely to be uniform. Ensuring fairness and transparency is a challenge that will need to be addressed.

## References

- [1] Ziqing Zhu, Ze Hu, Ka Wing Chan, Siqi Bu, Bin Zhou, and Shiwei Xia. Reinforcement learning in deregulated energy market: A comprehensive review. *Applied Energy*, 345:120360, 2023.
- [2] ATD Perera and Parameswaran Kamalaruban. Applications of reinforcement learning in energy systems. *Renewable and Sustainable Energy Reviews*, 137:110618, 2021.
- [3] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [4] Stefano V. Albrecht, Filippos Christianos, and Lukas Schäfer. *Multi-Agent Reinforcement Learning: Foundations and Modern Approaches*. MIT Press, 2024.
- [5] Lloyd S Shapley. Stochastic games. *Proceedings of the national academy of sciences*, 1953.
- [6] Nick Harder, Anke Weidlich, and Philipp Staudt. Finding individual strategies for storage units in electricity market models using deep reinforcement learning. *Energy Informatics*, 6(Suppl 1):41, 2023.
- [7] Thomas Wolgast and Astrid Nieße. Approximating energy market clearing and bidding with model-based reinforcement learning. *arXiv preprint arXiv:2303.01772*, 2023.
- [8] Panagiotis Michailidis, Iakovos Michailidis, and Elias Kosmatopoulos. Reinforcement Learning for Optimizing Renewable Energy Utilization in Buildings: A Review on Applications and Innovations. *Energies*, 18(7):1724, 2025.
- [9] Victor Ahlqvist, Pär Holmberg, and Thomas Tangerås. A survey comparing centralized and decentralized electricity markets. *Energy Strategy Reviews*, 40:100812, 2022.
- [10] Lucien Werner and Peeyush Kumar. Multi-market energy optimization with renewables via reinforcement learning. *arXiv preprint arXiv:2306.08147*, 2023.
- [11] Chenxiao Guan, Yanzhi Wang, Xue Lin, Shahin Nazarian, and Massoud Pedram. Reinforcement learning-based control of residential energy storage systems for electric bill minimization. In *2015 12th Annual IEEE Consumer Communications and Networking Conference (CCNC)*, pages 637–642. IEEE, 2015.
- [12] José R Vázquez-Canteli and Zoltán Nagy. Reinforcement learning for demand response: A review of algorithms and modeling techniques. *Applied energy*, 235:1072–1089, 2019.
- [13] Xin Qiu, Tu A Nguyen, and Mariesa L Crow. Heterogeneous energy storage optimization for microgrids. *IEEE Transactions on Smart Grid*, 7(3):1453–1461, 2015.
- [14] Elinor Ginzburg-Ganz, Itay Segev, Alexander Balabanov, Elior Segev, Sivan Kaully Naveh, Ram Machlev, Juri Belikov, Liran Katzir, Sarah Keren, and Yoash Levron. Reinforcement learning model-based and model-free paradigms for optimal control problems in power systems: Comprehensive review and future directions. *Energies*, 17(21):5307, 2024.
- [15] Ting Yang, Liyuan Zhao, Wei Li, and Albert Y Zomaya. Dynamic energy dispatch strategy for integrated energy system based on improved deep reinforcement learning. *Energy*, 235:121377, 2021.
- [16] Bin Zhang, Weihao Hu, Di Cao, Qi Huang, Zhe Chen, and Frede Blaabjerg. Deep reinforcement learning-based approach for optimizing energy conversion in integrated electrical and heating system with renewable energy. *Energy conversion and management*, 202:112199, 2019.
- [17] Aviad Navon, Juri Belikov, Ariel Orda, and Yoash Levron. On the stability of strategic energy storage operation in wholesale electricity markets. *arXiv preprint arXiv:2402.02428*, 2024.
- [18] Flora Charbonnier, Thomas Morstyn, and Malcolm D McCulloch. Coordination of resources at the edge of the electricity grid: Systematic review and taxonomy. *Applied Energy*, 318:119188, 2022.

- [19] Dimitrios Papadaskalopoulos and Goran Strbac. Nonlinear and randomized pricing for distributed management of flexible loads. *IEEE Transactions on Smart Grid*, 7(2):1137–1146, 2015.
- [20] Aisling Pigott, Constance Crozier, Kyri Baker, and Zoltan Nagy. Gridlearn: Multiagent reinforcement learning for grid-aware building energy management. *Electric Power Systems Research*, 2022.
- [21] Takao et al. Moriyama. Reinforcement learning testbed for power-consumption optimization. In *Methods and Applications for Modeling and Simulation of Complex Systems: 18th Asia Simulation Conference (AsiaSim)*. Springer, 2018.
- [22] José R. Vázquez-Canteli, Jérôme Kämpf, Gregor Henze, and Zoltan Nagy. Citylearn v1.0: An openai gym environment for demand response with deep reinforcement learning. Association for Computing Machinery, 2019.
- [23] Antoine et al. Marot. Learning to run a power network challenge: a retrospective analysis. In *NeurIPS 2020 Competition and Demonstration Track*, 2021.
- [24] Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U Balis, Gianluca De Cola, Tristan Deleu, Manuel Goulão, Andreas Kallinteris, Markus Krimmel, Arjun KG, et al. Gymnasium: A standard interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032*, 1:1–8, 2024.

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: Both sections promise (i) a formal hybrid-market model, (ii) a multi-agent RL framework for the ISO and strategic players, (iii) the new, modular Energy-Net simulator, and (iv) empirical validation across deterministic, stochastic, and strategic settings; these are delivered in Sections 3–7, so the claims match the actual contributions.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The discussion in the results section and the conclusion section discusses the limitation of our work. In addition, simplifying assumptions are mentioned throughout the paper.

### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The manuscript is empirical and algorithmic in focus, it contains no formal theorems or proofs, only modelling equations and simulation results, so the checklist item on theoretical assumptions and proofs is not applicable.

### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Section 7 the empirical evaluation specifies every simulator, configuration, hyperparameters, and algorithm setting. We also provide our code base and full result in the appendix.

### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The code and configuration files as well as the full results are in the supplementary materials. In addition, we will release an open-access repository (with the `run.sh` reproducer and all YAML scenario files) alongside the camera-ready version.

### 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Section 7 list all hyper-parameters, network/optimizer choices, and scenario settings.

### 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: All settings were run multiple times to achieve statistical significance, for which the values are reported.

#### 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

We specify our resources in Section 7.

#### 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The work relies solely on synthetic simulations, no personal or proprietary data are used, no human subjects are involved, and all methods pose no foreseeable societal or environmental harm

#### 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: Section 8 outlines the positive potential and flags possible negatives

#### 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: The paper releases only a power-system simulation environment and illustrative RL code; no large pretrained models or scraped datasets with dual-use or safety concerns are provided, so additional safeguards are not applicable.

#### 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All third-party assets are cited with their official references and license terms

#### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: The paper releases a new simulator full described in the paper (code is in the supplementary material). The GitHub repository includes: a README with installation and run commands, API docs for every module, default configuration files. No personal data are involved, so no consent was required.

#### 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The study involves only simulated power-system agents; no human participants or crowdsourced data were used.

**15. Institutional review board (IRB) approvals or equivalent for research with human subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The research relies exclusively on simulated agents and publicly available data; no experiments involving human subjects were conducted, so IRB approval is not applicable.

**16. Declaration of LLM usage**

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: Large Language Models were not part of the proposed market-optimization framework; any LLM use was limited to routine writing assistance and had no bearing on the methodology or results.

## A RL and MARL

A Reinforcement Learning (RL) problem can be defined as a Markov Decision Process (MDP) represented by the tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$ , where:

- $\mathcal{S}$  is the set of states,
- $\mathcal{A}$  is the set of actions,
- $\mathcal{P}(s' | s, a)$  is the transition probability from state  $s$  to  $s'$  under action  $a$ ,
- $\mathcal{R}(s, a)$  is the reward function,
- $\gamma \in [0, 1]$  is the discount factor.

The goal is to find a policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  that maximizes the expected cumulative reward:

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^t R_t \right],$$

where  $R_t$  is the reward received at time step  $t$ . It is assumed that the MDP is too large to efficiently compute  $\pi^*$ , so approximation methods are employed to estimate it. These methods often involve learning value functions or directly optimizing parameterized policies using sampled interactions with the environment.

The problem can be modeled as a Markov Decision Process (MDP), defined by the tuple:

$$\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$$

where:

- $\mathcal{S}$ : The set of states, defined by  $\mathcal{S} = \{(t, \sigma_t) | t = 1, \dots, T, 0 \leq \sigma_t \leq S_{\max}\}$ ,
- $\mathcal{A}$ : The set of actions, where each action is represented by the pair  $(P_t^b, P_t^s)$ ,
- $\mathcal{P}(s' | s, a)$ : The state transition function, given by:

$$\mathcal{P}(s' | s, a) = \Pr(\sigma_{t+1} | \sigma_t, P_t^b, P_t^s),$$

- $\mathcal{R}(s, a)$ : The reward function:

$$\mathcal{R}(s, a) = \phi_t(P_t^s) - \xi_t(P_t^b),$$

- $\gamma$ : The discount factor,  $\gamma \in [0, 1]$ , which determines the relative importance of future rewards.

The goal is to find an optimal policy  $\pi^*$  that maximizes the expected cumulative reward:

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=1}^T \gamma^{t-1} \mathcal{R}(s_t, a_t) \right],$$

where:

- $s_t = (t, \sigma_t)$  is the state at time  $t$ ,
- $a_t = (P_t^b, P_t^s)$  is the action at time  $t$ ,
- $\mathcal{R}(s_t, a_t)$  is the immediate reward obtained from taking action  $a_t$  in state  $s_t$ .

rl and marl algorithms can be broadly categorized as model-free, which learn policies directly from experience without modeling the environment, and model-based, which learn or use environment models to plan or simulate outcomes. Model-free methods (e.g., value-based or policy gradient) tend to be more scalable but sample-inefficient, while model-based methods improve sample efficiency and enable planning but struggle with modeling complex dynamics [4].

Reinforcement Learning (rl) is a learning paradigm where an agent learns optimal behavior by interacting with an environment and receiving rewards or penalties for its actions [3] (see Appendix ?? for a full definition). Multi-agent RL (marl) extends rl to scenarios involving multiple autonomous agents that concurrently learn and make decisions within a shared or partially shared environment. Each agent aims to maximize its own utility (typically measured as accumulated reward), but its actions can influence both its own outcomes and the outcomes of other agents, leading to complex emergent behaviors and the need for coordination and cooperation.

marl is particularly suitable for modeling energy systems and networks, since they are inherently multi-agent environments composed of diverse, distributed, and strategically autonomous entities, such as grid-edge components and prosumers, utility companies, system operators, and market participants. These entities have different objectives, interact over shared physical and economic infrastructures, and must respond dynamically to system conditions, prices, and regulations. MARL provides a natural framework to model these interactions, enabling agents to learn adaptive policies, coordinate under uncertainty, and reason about both cooperative and competitive dynamics. Moreover, its ability to simulate emergent behavior and explore decentralized strategies makes it a powerful tool for both designing and analyzing modern energy systems.

The most common marl model is the Stochastic Game (sg) (also known as Markov Game or Multi-agent MDP) [5] defined as a tuple  $\langle \mathcal{S}, \mathcal{A} = \{\mathcal{A}_i\}_{i=1}^n, \mathcal{T}, \mathcal{R} = \{\mathcal{R}_i\}_{i=1}^n, \gamma \rangle$ , where  $\mathcal{S}$  is the state space,  $\mathcal{A}$  is the *joint action space* with  $\mathcal{A}_i$  as the  $i^{th}$  agent action space s.t.  $a \triangleq (a_1, a_2, \dots, a_n)$  for  $a \in \mathcal{A}$ ,  $\mathcal{T} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  is the transition probability function  $\mathcal{T}(s', a, s)$  such that  $\forall s \in \mathcal{S}, \forall a \in \mathcal{A} : \sum_{s' \in \mathcal{S}} \mathcal{T}(s, a, s') = 1$ ,  $\mathcal{R}$  is the *joint reward function* with  $\mathcal{R}_i : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  as the  $i^{th}$  agent reward function, and  $\gamma \in [0, 1]$  is the discount factor. A solution is a joint policy  $\pi \triangleq (\pi_1, \dots, \pi_n)$  associating each agent with policy  $\pi_i : \mathcal{S} \times \mathcal{A}_i \rightarrow [0, 1]$  that specifies the probability of agent  $i$  taking an action at a given state. The joint policy should achieve certain conditions on the expected returns yielded to agents (e.g., Nash equilibrium) [4]. The value (utility) function  $V_i^\pi(s)$  denotes the expected cumulative discounted reward agent  $i$  receives when starting in state  $s$  and the agents follow joint policy  $\pi$  thereafter. The action-value function or Q-value  $Q_i^\pi(s, a)$  extends this notion by quantifying the expected value when performing  $a$  in  $s$ , and then continuing according to  $\pi$ . A Multi-agent Partially Observed MDP (or Partially Observable Stochastic Game) also includes for each agent observation set  $O_i$  and a sensor function  $\mathcal{O}_i : \mathcal{A} \times \mathcal{S} \times O_i \rightarrow [0, 1]$ .

This general definition captures a variety of interactions and relationships that can exist between agents in collaborative, competitive, and mixed-incentive MARL settings. Complex agent interactions may give rise to behaviors that are difficult to anticipate by simply examining each agent in isolation. Thus, despite the potential to solve complex problems across various domains, marl faces various significant challenges that stem from aspects such as scale, conflicting goals of self-interested agents, and the concurrent learning of the different agents [4]. All these are relevant to MARL in general but are particularly relevant to energy networks with the added need to account for the dynamics of the physical environment and the effect decisions may have on the functioning of the electricity network.

29 RL and MARL algorithms can be broadly categorized as model-free, which learn policies directly 130 from experience without modeling the environment, and model-based, which learn or use environment 131 models to plan or simulate outcomes. Model-free methods (e.g., value-based or policy gradient) tend 132 to be more scalable but sample-inefficient, while model-based methods improve sample efficiency 133 and enable planning but struggle with modeling complex dynamics [? ]

## B Energy Market Dynamics

### B.1 Energy Markets and the Dispatch Problem

Historically, the energy market comprised three principal components: power producers (e.g., power plants), power consumers (industrial and residential), and the ISO, responsible for market management and coordination. The producers typically used conventional coal-based generation and were either units under the full control of the ISO, or independent units that participated in the market but that were fully regulated, i.e., bound by production agreements made with the .

A typical structure of a market was based on the *day-ahead market* in which the ISO predicts the following day's power demand and issues a *dispatch*, an offline production schedule to each producer while considering operational constraints and generation costs. The dispatch traditionally divides the 24-hour planning horizon into 48 discrete half-hour time periods. In addition to the generation of the predicted, or *nominal* demand, the ISO also manages the *reserve*, which sets a backup production capability for each time step. If in real-time the

controlled production determined by the dispatch is not enough to cover the realized demand, reserves, which are more flexible but also more expensive and polluting, are activated by an online controller. Producers are then compensated based on the System Marginal Price (SMP) mechanism, which is calculated as the marginal cost of producing the final unit of energy required to satisfy system demand, based on the least-cost dispatch solution (See Appendix ??). For the purposes of this work we abstract the dispatch details, and consider only the total amount of power produced at each timestamp, as well as its total cost to the ISO with no regard to the inner structure of the dispatch.

Recent reforms in the power market have introduced independent grid-edge market players, which we denote as **GEAgents**, including private electric companies and smart homes. These new market players possess the ability to produce electricity, manage internal consumption, and utilize power storage capabilities. Unlike traditional controlled producers, they are not legally required to adhere to dispatch instructions and may buy from or sell to the grid at will to maximizing their profits. We assume GEAgents are rational, so the natural way for the ISO to induce desired behaviors of the market players is via price signals. In real-time operations, the ISO manages the grid by buying electricity from power producers and selling it to consumers. The selling price at time  $t$ , denoted as  $\xi_t$ , and the feed-in price, denoted as  $\phi_t$ , are the primary tools for market control.

The GEAgent models are essential for the ISO's planning, as they capture participant strategies and behaviors that influence the grid's supply-demand balance. These models enable the ISO to design pricing mechanisms, such as sell prices and feed-in tariffs, to align player incentives with grid stability and efficiency. We classify market player behaviors in increasingly realistic environments, starting with simpler cases to build intuition before progressing to more complex scenarios, as the problems share similar structures. In correspondence with current energy markets, each GEAgent operates a Production-Consumption-Storage (PCS) unit, which can produce (e.g., via pv), consume (e.g., via electrical appliances) and store (e.g., via a battery) energy. It aims at maximizing its profit over the period in question.

To determine dispatch and pricing, the ISO utilizes demand predictions for the subsequent 24 hours, denoted  $\hat{D}_t$ , where  $t$  represents the time interval. Based on these predictions, the ISO determines a scheduled production dispatch  $\Delta_t$  for each timestamp. It also determines for each time step how much reserve to guarantee, specifying the standby capacity to maintain in response to unexpected demand surges or generation outages. Reserve energy enhances grid reliability but can be highly polluting when supplied by fossil-fuel generators, which operate inefficiently and emit more greenhouse gases.

In real-time operations, the ISO manages the grid by buying electricity from power producers and selling it to consumers. The selling price at time  $t$ , denoted as  $\xi_t$ , and the feed-in price, denoted as  $\phi_t$ , are the primary tools for market control.

The electricity market includes  $n$  independent agents representing the GEAgents, indexed by  $i \in \{1, \dots, n\}$ , who operate autonomously to maximize their profits. The ISO has no direct control over these agents, and their interactions are governed by market dynamics, which are influenced by various regulations. These regulations, coupled with non-economic factors, significantly shape the cost structure of the system. However, the ISO can compute costs based on relevant inputs and adapt its computational models dynamically to reflect changes in regulations or legislation.

In this work, we suggest using RL-to model the market participants and ways for the ISO to control the dispatch  $\Delta$  and price signals  $\xi, \phi$  to minimize the total costs for the ISO(thus the taxpayers) while satisfying the demand. A key challenge is that this needs to be done while taking market players' strategic behavior into account.

To support optimizing the ISO's behavior, we analyze how market players react to prices in increasingly complex settings, from deterministic to stochastic and strategic environments.

## B.2 Market Participants

The GEAgent models are essential for the ISO's planning, as they capture participant strategies and behaviors that influence the grid's supply-demand balance. These models enable the ISO to design pricing mechanisms, such as sell prices and feed-in tariffs, to align player incentives with grid stability and efficiency. We classify market player behaviors in increasingly realistic environments, starting with simpler cases to build intuition before progressing

to more complex scenarios, as the problems share similar structures. In correspondence with current energy markets, each GEAgent operates a Production-Consumption-Storage (PCS) unit, which can produce (e.g., via pv), consume (e.g., via electrical appliances) and store (e.g., via a battery) energy. It aims at maximizing its profit over the period in question.

Having settled on market players', we proceed to present the task that the ISO faces. The ISO is tasked with meeting electricity demand at all times. To achieve this, the ISO controls the dispatch of electricity generation. While the specifics of which power plant generates how much power are abstracted, the total scheduled electricity production is determined for each time step, ensuring sufficient supply to meet anticipated demand.

The ISO aims to maximize its utility, which may include balancing grid supply and demand, minimizing operational costs, or promoting renewable integration.

The total cost incurred by the ISO increases marginally due to the characteristics of the SMP mechanism. The SMP prioritizes electricity from the cheapest sources first, resulting in higher costs for additional megawatts of production as cheaper resources are exhausted. Additionally, sharp changes in production across time steps introduce significant costs due to ramp-up and cool-down constraints of power plants. These transitions strain generation units, necessitating increased operational expenses. The ISO incorporates these costs into pricing to discourage abrupt fluctuations, maintaining grid stability.

To influence the behavior of market players, the ISO offers sell prices and feed-in tariffs. These prices act as economic signals, encouraging players to adjust their electricity consumption, production, and storage behaviors in alignment with grid stability and efficiency goals. By strategically setting these prices, the ISO aims to optimize the overall operation of the electricity market under a hybrid public-private model.

In what follows, we examine three levels of complexity that are associated with the nature and pattern of the demand (consumption): deterministic and known, stochastic, and strategic.

### B.3 Deterministic Setting

As a first step, we consider a fully deterministic environment, where the demand is fully known in advance and the prices are set in advance (at time 0 of every day).

- Storage capacity:  $S_{\max}$ .
- Maximum charging rate:  $C_{\max}$ .
- Maximum discharging rate:  $D_{\max}$ .
- Initial storage state of charge:  $\sigma_0$ .
- Selling price levels  $\xi_t$  set by the ISO for each time interval and known in advance to the player.
- Feed-in prices  $\phi_t$  set by the ISO and known in advance to the player as well.

Since all information is given in advance, the GEAgent can compute optimal policies at time-step 0. A GEAgent must decide how much power to buy from ( $P_t^b$ ), and sell to ( $P_t^s$ ) the grid at every timestamp  $t$  to maximize its total revenue under its operational constraints. Formally:

$$\max \sum_{t=1}^T (\phi_t(P_t^s) - \xi_t(P_t^b)) \quad (\text{Deterministic GEAgent Objective})$$

Subject to:

(a) **Power Balance Constraints:**

At each time  $t$ , the power bought or sold must meet the demand, including charging:

$$\forall t : P_t^b - P_t^s = l_t + (\sigma_t - \sigma_{t-1}) \quad (C1)$$

(b) **Storage Capacity Constraints:**

The storage level must remain within capacity limits:

$$\forall t : 0 \leq \sigma_t \leq S_{\max} \quad (C2)$$

(c) **Charging and Discharging Rate Constraints:**

$$\forall t : -D_{\max} \leq P_t^b - (l_t + P_t^s) \leq C_{\max} \quad (C3)$$

(d) **Non-Negativity Constraints:**

$$\forall t : P_t^b, P_t^s, \sigma_t \geq 0 \quad (\text{C4})$$

(e) **No Simultaneous Charging and Discharging:**

$$\forall t : P_t^b \cdot P_t^s = 0 \quad (\text{C5})$$

**ISO** In the deterministic case, at the start of the planning horizon (timestamp 0), the ISO receives the following inputs:

- Demand  $D_t$  for all timestamps in the horizon.
- Reserve activation cost  $C_{\text{reserve}}$ .
- The number of GEAgents  $N$  participating in the market.
- The maximum discharge rates of each market player  $D_{\max}^i$ .

Based on this information, the ISO determines the scheduled amount of production  $\Delta_t$  and prices  $\xi_t(\cdot)$ ,  $\phi_t(\cdot)$  for all timestamps  $t \in [T]$  ahead. Then, at each timestamp  $t$  market players can respond to the prices by buying or selling power to the grid, contributing a net power demand  $P_t^{\text{net}}$ . If the net demand after accounting for  $P_t^{\text{net}}$  exceeds the scheduled production  $\Delta_t$ , the ISO activates reserves or peaker plants to cover the shortfall. If the market players are assumed to be rational, and the ISO makes the prices public at  $t = 0$ , the market players are solving the deterministic problem as presented in Section B.3, and the ISO can run the simulation of the market players to optimize the dispatch and the price signal. An example of this behavior and its benefits is given in Section ??.

The ISO aims to minimize its total costs,

$$\min C^{\text{total}} = \min \left[ C^{\text{dispatch}} + \sum_{t=1}^T C_t^{\text{online}} \right] \quad (\text{ISO objective})$$

where:

- **Cost of the Dispatch Schedule ( $C^{\text{dispatch}}$ ):**

$$C^{\text{dispatch}} = \sum_{t=1}^T C(\Delta_t) + \sum_{t=2}^T \rho(\Delta_0, \dots, \Delta_t),$$

where  $\rho$  is a penalty function that can be tailored to various performance criteria, e.g., for penalizing sharp changes in dispatch levels between consecutive periods.

- **Online Cost per Timeframe ( $C_t^{\text{online}}$ ):** The sum of the market cost and the reserve activation cost:

$$C_t^{\text{online}} = C_t^{\text{market}} + C_t^{\text{reserve}}(\max(0, D_t - P_t^{\text{net}} - \Delta_t)),$$

Notably, we assume that all demand must be met, a constraint that can be relaxed if needed.

- **Market Cost per Timeframe ( $C_t^{\text{market}}$ ):** Payments to market players for the power they sell to the grid net of the revenue from selling the power to market players:

$$C_t^{\text{market}} = \sum_i \phi_t^{(i)}(s_t^{(i)}) - \sum_i \xi_t^{(i)}(b_t^{(i)})$$

where  $\phi_t^{(i)}$  is the feed-in tariff offered to player  $i$  at time  $t$ , and  $s_t^{(i)}$  is the amount of power sold by player  $i$  to the grid.

Note that this problem is unconstrained, since we assume that when the demand is not met by the production and the market, the ISO operates the reserves. The incentive to meet the demand using generation is encapsulated in the typically high costs associated with activating the reserves.

## B.4 Accounting for Stochasticity

Real-world systems are inherently stochastic, requiring models to account for uncertainty. Key sources of randomness include:

- Internal load variability,
- Renewable production fluctuations,
- Price changes driven by external demand uncertainty.

All these may lead to an inability to exactly predict the demand that will be needed.

From the point of view of the GEAgent, the main source of uncertainty can come from its To address this, the objective function is reformulated as:

$$\max \mathbb{E}_{l_t, \xi_t} \left[ \sum_{t=1}^T (\phi_t(P_t^s) - \xi_t(P_t^b)) \right]. \quad (\text{Stochastic Player Objective})$$

At each timestamp  $t$ , the player observes the realizations of  $l_t$ ,  $g_t$ , and  $\xi_t$  before deciding on  $P_t^b$ , and  $P_t^s$ .

In a stochastic environment, the distributions of  $l_t$ ,  $g_t$ , and  $\xi_t$  may be unknown. If this is the case, the player can estimate these distributions from historical data and observations using machine learning methods to improve decision-making under these forms of uncertainty.

The main change becomes the uncertainty about the demand

In this case, we have two options, depending on when decisions need to be made.

## B.5 Accounting for Load Flexibility and Strategic Demand

So far, we considered settings in which all participants were aiming to maximize their revenue (and minimize cost) while considering the deterministic or stochastic information that is received at time step 0, i.e., at the beginning of the daily episode. This meant that prices and dispatch decisions are made at the start of each episode, with the real-time decisions limited to reserve activation or curtailment (energy discharge) actions in response to unpredictable demand and the requirement to maintain stability.

In modern energy systems, demand is not only stochastic but also strategic. This is because grid-edge agents can intelligently manage the operation of devices and distributed energy resources (DERs), in response to system-level signals, such as prices, frequency, or voltage. This *load flexibility* is reshaping energy markets by introducing new ways by which grid-edge agents can contribute to the efficient and stable operation of the network [18, 1]. However, this shift also introduces challenges such as increased system complexity, uncertainty in demand forecasting, and the need for regulatory mechanisms to ensure fair and reliable participation.

In this extended setting, the ISO aims to maximize its utility, but needs to determine the selling price  $\xi_t$  and feed-in prices  $\phi_t$  for each time step  $t$  according to the demand  $D_t$  at time  $t$ . The key challenge is that  $D_t$  now includes the GEAagents ability to sell, buy, and store power. From the perspective of the GEAagent, the price signals  $\xi_t(P_t^s, P_t^b, \dots)$  represent the exogenous prices set by the ISO, which depend on the player's sales  $P_t^s$  and purchases  $P_t^b$  as well as other variables. This coupling results in a feedback mechanism where the player's actions influence the prices, and the prices in turn affect the player's actions. This introduces a game-theoretic dimension to the problem that the market player faces, where the player's decisions on  $P_t^b$ , and  $P_t^s$  are influenced by the GSO's pricing strategy and vice versa.

It is important to clarify what the possibilities are that are available to the ISO with regard to the dispatch and pricing decisions it can make. This is not only a technical question, but a regulatory and policy-making question that needs to be accounted for. Two common approaches are day-ahead and dynamic pricing.

Formally, the GEAagent's input includes timesteps  $t = 1, 2, \dots, T$  GEAagent's load:  $l_t$ , storage capacity:  $S_{\max}$ , maximum charging rate:  $C_{\max}$ , maximum discharging rate:  $D_{\max}$ , current storage state of charge:  $\sigma_0$  as defined in sections B.3 and B.4. The key difference is that now the selling price  $\xi_t$  and feed-in prices  $\phi_t$  can be set by ISOin advance or in a dynamic way, in response to the market state.

The objective is now:

$$\max_{P_t^b, P_t^s} \mathbb{E}_{l_t, g_t} \left[ \sum_{t=1}^T (\phi_t(P_t^s, P_t^b, \dots) - \xi_t(P_t^s, P_t^b, \dots)) \right]. \quad (\text{Strategic Player Objective})$$

The ISO at the start of the planning horizon (timestamp 0), the ISO receives the following inputs:

- The cost function of the production  $C(\Delta_t)$  for each  $t$ .
- Predicted demand  $\hat{D}_t$  for all timestamps in the horizon.
- Reserve activation cost per unit  $C_{\text{reserve}}$ .
- The number of market players  $N$  participating in the market.
- The maximum discharge rates of each market player  $D_{\max}^i$ .

Based on this information, the ISO determines the scheduled amount of production  $\Delta_t$  for each timestamp in the horizon. Here, it is crucial to distinguish between nominal and flexible demand components. Nominal demand, denoted  $D$  refers to the exogenous, inelastic portion of load at each grid node that remains unaffected by local control strategies, real-time market incentives, or variations in renewable generation. In contrast, *flexible demand*, denoted  $l$ , refers to the portion of demand (electricity consumption) that can be adjusted in time, quantity, or pattern in response to external signals—such as price changes, grid conditions, or availability of renewable energy.

The objective of the ISO now becomes

$$\min \mathbb{E}_{D, l} \left[ C^{\text{dispatch}} + \sum_{t=1}^T C_t^{\text{online}} \right] \quad (\text{O2})$$

Since it is impossible for the ISO to precisely model market players' demand without considering its strategic nature, optimization methods that are appropriate for deterministic and stochastic settings won't work here. Thus, as we specify in the next section, we model the market using RL.

## B.6 Deterministic Setting

As a first step, we consider a fully deterministic environment, where the demand is fully known in advance and the prices are set in advance (at time 0 of every day).

- Storage capacity:  $S_{\max}$ .
- Maximum charging rate:  $C_{\max}$ .
- Maximum discharging rate:  $D_{\max}$ .
- Initial storage state of charge:  $\sigma_0$ .
- Selling price levels  $\xi_t$  set by the ISO for each time interval and known in advance to the player.
- Feed-in prices  $\phi_t$  set by the ISO and known in advance to the player as well.

Since all information is given in advance, the GEAgent can compute optimal policies at time-step 0. A GEAgent must decide how much power to buy from  $(P_t^b)$ , and sell to  $(P_t^s)$  the grid at every timestamp  $t$  to maximize its total revenue under its operational constraints. Formally:

$$\max \sum_{t=1}^T (\phi_t(P_t^s) - \xi_t(P_t^b)) \quad (\text{Deterministic GEAgent Objective})$$

Subject to:

- (a) **Power Balance Constraints:**

At each time  $t$ , the power bought or sold must meet the demand, including charging:

$$\forall t : P_t^b - P_t^s = l_t + (\sigma_t - \sigma_{t-1}) \quad (\text{C1})$$

(b) **Storage Capacity Constraints:**

The storage level must remain within capacity limits:

$$\forall t : 0 \leq \sigma_t \leq S_{\max} \quad (\text{C2})$$

(c) **Charging and Discharging Rate Constraints:**

$$\forall t : -D_{\max} \leq P_t^b - (l_t + P_t^s) \leq C_{\max} \quad (\text{C3})$$

(d) **Non-Negativity Constraints:**

$$\forall t : P_t^b, P_t^s, \sigma_t \geq 0 \quad (\text{C4})$$

(e) **No Simultaneous Charging and Discharging:**

$$\forall t : P_t^b \cdot P_t^s = 0 \quad (\text{C5})$$

**ISO** In the deterministic case, at the start of the planning horizon (timestamp 0), the ISO receives the following inputs:

- Demand  $D_t$  for all timestamps in the horizon.
- Reserve activation cost  $C_{\text{reserve}}$ .
- The number of GEAgents  $N$  participating in the market.
- The maximum discharge rates of each market player  $D_{\max}^i$ .

Based on this information, the ISO determines the scheduled amount of production  $\Delta_t$  and prices  $\xi_t(\cdot)$ ,  $\phi_t(\cdot)$  for all timestamps  $t \in [T]$  ahead. Then, at each timestamp  $t$  market players can respond to the prices by buying or selling power to the grid, contributing a net power demand  $P_t^{\text{net}}$ . If the net demand after accounting for  $P_t^{\text{net}}$  exceeds the scheduled production  $\Delta_t$ , the ISO activates reserves or peaker plants to cover the shortfall. If the market players are assumed to be rational, and the ISO makes the prices public at  $t = 0$ , the market players are solving the deterministic problem as presented in Section B.3, and the ISO can run the simulation of the market players to optimize the dispatch and the price signal. An example of this behavior and its benefits is given in Section ??.

The ISO aims to minimize its total costs,

$$\min C^{\text{total}} = \min \left[ C^{\text{dispatch}} + \sum_{t=1}^T C_t^{\text{online}} \right] \quad (\text{ISO objective})$$

where:

- **Cost of the Dispatch Schedule ( $C^{\text{dispatch}}$ ):**

$$C^{\text{dispatch}} = \sum_{t=1}^T C(\Delta_t) + \sum_{t=2}^T \rho(\Delta_0, \dots, \Delta_t),$$

where  $\rho$  is a penalty function that can be tailored to various performance criteria, e.g., for penalizing sharp changes in dispatch levels between consecutive periods.

- **Online Cost per Timeframe ( $C_t^{\text{online}}$ ):** The sum of the market cost and the reserve activation cost:

$$C_t^{\text{online}} = C_t^{\text{market}} + C_t^{\text{reserve}}(\max(0, D_t - P_t^{\text{net}} - \Delta_t)),$$

Notably, we assume that all demand must be met, a constraint that can be relaxed if needed.

- **Market Cost per Timeframe ( $C_t^{\text{market}}$ ):** Payments to market players for the power they sell to the grid net of the revenue from selling the power to market players:

$$C_t^{\text{market}} = \sum_i \phi_t^{(i)}(s_t^{(i)}) - \sum_i \xi_t^{(i)}(b_t^{(i)})$$

where  $\phi_t^{(i)}$  is the feed-in tariff offered to player  $i$  at time  $t$ , and  $s_t^{(i)}$  is the amount of power sold by player  $i$  to the grid.

Note that this problem is unconstrained, since we assume that when the demand is not met by the production and the market, the ISO operates the reserves. The incentive to meet the demand using generation is encapsulated in the typically high costs associated with activating the reserves.

## B.7 Accounting for Stochasticity

Real-world systems are inherently stochastic, requiring models to account for uncertainty. Key sources of randomness include:

- Internal load variability,
- Renewable production fluctuations,
- Price changes driven by external demand uncertainty.

All these may lead to an inability to exactly predict the demand that will be needed.

From the point of view of the GEAgent, the main source of uncertainty can come from its To address this, the objective function is reformulated as:

$$\max \mathbb{E}_{l_t, \xi_t} \left[ \sum_{t=1}^T (\phi_t(P_t^s) - \xi_t(P_t^b)) \right]. \quad (\text{Stochastic Player Objective})$$

At each timestamp  $t$ , the player observes the realizations of  $l_t$ ,  $g_t$ , and  $\xi_t$  before deciding on  $P_t^b$ , and  $P_t^s$ .

In a stochastic environment, the distributions of  $l_t$ ,  $g_t$ , and  $\xi_t$  may be unknown. If this is the case, the player can estimate these distributions from historical data and observations using machine learning methods to improve decision-making under these forms of uncertainty.

The main change becomes the uncertainty about the demand

In this case, we have two options, depending on when decisions need to be made.

## B.8 Accounting for Load Flexibility and Strategic Demand

So far, we considered settings in which all participants were aiming to maximize their revenue (and minimize cost) while considering the deterministic or stochastic information that is received at time step 0, i.e., at the beginning of the daily episode. This meant that prices and dispatch decisions are made at the start of each episode, with the real-time decisions limited to reserve activation or curtailment (energy discharge) actions in response to unpredictable demand and the requirement to maintain stability.

In modern energy systems, demand is not only stochastic but also strategic. This is because grid-edge agents can intelligently manage the operation of devices and distributed energy resources (DERs), in response to system-level signals, such as prices, frequency, or voltage. This *load flexibility* is reshaping energy markets by introducing new ways by which grid-edge agents can contribute to the efficient and stable operation of the network [18, 1]. However, this shift also introduces challenges such as increased system complexity, uncertainty in demand forecasting, and the need for regulatory mechanisms to ensure fair and reliable participation.

In this extended setting, the ISO aims to maximize its utility, but needs to determine the selling price  $\xi_t$  and feed-in prices  $\phi_t$  for each time step  $t$  according to the demand  $D_t$  at time  $t$ . The key challenge is that  $D_t$  now includes the GEAagents ability to sell, buy, and store power. From the perspective of the GEAagent, the price signals  $\xi_t(P_t^s, P_t^b, \dots)$  represent the exogenous prices set by the ISO, which depend on the player's sales  $P_t^s$  and purchases  $P_t^b$  as well as other variables. This coupling results in a feedback mechanism where the player's actions influence the prices, and the prices in turn affect the player's actions. This introduces a game-theoretic dimension to the problem that the market player faces, where the player's decisions on  $P_t^b$ , and  $P_t^s$  are influenced by the GSO's pricing strategy and vice versa.

It is important to clarify what the possibilities are that are available to the ISO with regard to the dispatch and pricing decisions it can make. This is not only a technical question, but a regulatory and policy-making question that needs to be accounted for. Two common approaches are day-ahead and dynamic pricing.

Formally, the GEAagent's input includes timesteps  $t = 1, 2, \dots, T$  GEAagent's load:  $l_t$ , storage capacity:  $S_{\max}$ , maximum charging rate:  $C_{\max}$ , maximum discharging rate:  $D_{\max}$ , current storage state of charge:  $\sigma_0$  as defined in sections B.3 and B.4. The key difference is that now the selling price  $\xi_t$  and feed-in prices  $\phi_t$  can be set by ISOin advance or in a dynamic way, in response to the market state.

The objective is now:

$$\max_{P_t^b, P_t^s} \mathbb{E}_{l_t, g_t} \left[ \sum_{t=1}^T (\phi_t(P_t^s, P_t^b, \dots) - \xi_t(P_t^s, P_t^b, \dots)) \right]. \quad (\text{Strategic Player Objective})$$

The ISO at the start of the planning horizon (timestamp 0), the ISO receives the following inputs:

- The cost function of the production  $C(\Delta_t)$  for each  $t$ .
- Predicted demand  $\hat{D}_t$  for all timestamps in the horizon.
- Reserve activation cost per unit  $C_{\text{reserve}}$ .
- The number of market players  $N$  participating in the market.
- The maximum discharge rates of each market player  $D_{\max}^i$ .

Based on this information, the ISO determines the scheduled amount of production  $\Delta_t$  for each timestamp in the horizon. Here, it is crucial to distinguish between nominal and flexible demand components. Nominal demand, denoted  $D$  refers to the exogenous, inelastic portion of load at each grid node that remains unaffected by local control strategies, real-time market incentives, or variations in renewable generation. In contrast, *flexible demand*, denoted  $l$ , refers to the portion of demand (electricity consumption) that can be adjusted in time, quantity, or pattern in response to external signals—such as price changes, grid conditions, or availability of renewable energy.

The objective of the ISO now becomes

$$\min \mathbb{E}_{D, l} \left[ C^{\text{dispatch}} + \sum_{t=1}^T C_t^{\text{online}} \right] \quad (\text{O2})$$

Since it is impossible for the ISO to precisely model market players' demand without considering its strategic nature, optimization methods that are appropriate for deterministic and stochastic settings won't work here. Thus, as we specify in the next section, we model the market using RL.

## C SMP

A typical structure of a market was based on the day-ahead market in which the ISO predicts the following day's power demand and issues a *dispatch*, an offline production schedule to each producer while considering operational constraints and generation costs. The dispatch traditionally divides the 24-hour planning horizon into 48 discrete half-hour time periods. In addition to the generation of the predicted or nominal demand, the ISO also manages the reserve, which sets a backup production capability for each time step. If in real-time the controlled production determined by the dispatch is not enough to cover the realized demand, reserves, which are more flexible but also more expensive and polluting, are activated by an online controller. Producers are then compensated based on the System Marginal Price (SMP) mechanism, which is calculated as the marginal cost of producing the final unit of energy required to satisfy system demand, based on the least-cost dispatch solution.

Formally, let:

- $P_t$  be the total power production at time  $t$ ,
- $D_t$  be the total system demand at time  $t$ ,
- $C(P_t)$  be the cost function for production.

The SMP at timestamp  $t$  is defined as:

$$\kappa_t = \frac{\partial C(P_t)}{\partial P_t} \Big|_{P_t=D_t},$$

where  $\kappa_t$  represents the marginal cost of meeting the demand  $D_t$  using the least-cost generation defined by the merit-order curve.

In electricity markets, the SMP clears the market by equating supply and demand while satisfying the economic dispatch problem:

$$\min_{P_t} C(P_t) \quad \text{subject to} \quad P_t = D_t.$$

The SMP ensures that all dispatched generators receive the same price, incentivizing efficiency and cost-reflective bidding in competitive electricity markets. Note that SMP is non-decreasing with respect to the amount of power produced, meaning higher power demand usually results in a higher price *per kWh*. Consequently, reducing peak consumption is critical for lowering overall costs in the electricity market.

## D Dynamic Programming Formulation for a Storage Only PCS-unit Agent

The dynamic programming formulation for the optimization problem for storage control is given as:

- **State Variables:**
  - Current time step  $t$ ,
  - Current storage level  $\sigma_t$ .
- **Decision Variables:**
  - Energy bought  $P_t^b$ ,
  - Energy sold  $P_t^s$ .
- **Transition Function:**

$$\sigma_{t+1} = \sigma_t + (P_t^b - l_t - P_t^s).$$
- **Objective Function:** The immediate reward at each time step is:

$$r(P_t^b, P_t^s) = \phi_t(P_t^s) - \xi_t(P_t^b).$$

The cumulative reward is maximized over all time steps.

- **Recurrence Relation:**

$$V(t, \sigma_t) = \max_{P_t^b, P_t^s} [r(P_t^b, P_t^s) + V(t+1, \sigma_{t+1})],$$

subject to the constraints.

Similar methods adapted for stochastic optimization could be employed for the case where distribution is either known or can be approximated from existing data. In the case of the stochastic demand, there may even be an ability to compute a contingent policy that would deal with the stochastic signals.

## E Quadratic Pricing

This example demonstrates the possible impact of price intervention on market dynamics. We assume deterministic setting for the ISO for clarity, but the same logic can be applied in the non-deterministic scenario. Drawing from the literature [19], we apply superlinear and sublinear pricing adjustments to selling and feed-in tariffs, respectively.

The selling price incorporates a superlinear component:

$$\xi_t = \lambda^{buy} * P_t^b + \beta * [P_t^b]^2,$$

where  $\lambda^{buy}$  is a baseline price. Similarly, the feed-in price adds a sublinear adjustment:

$$\phi_t = \lambda^{feedin} * P_t^s + \gamma * \sqrt{P_t^s},$$

where  $\lambda^{feedin}$  is the baseline feed-in price.

Once per episode, at  $t = 0$ , the ISO commits to six coefficients  $(\alpha_0, \alpha_1, \alpha_2, \beta_0, \beta_1, \beta_2)$  that instantiate the quadratic tariff  $\pi^{buy/sell}(x)$  in (??). These coefficients stay fixed for the ensuing  $T$  steps; dispatch tweaks  $\delta_t$  may still follow online if enabled.

**Baseline Scenario** Assume the demand structure as described by Table 2 and  $\rho = 0.3$ . Also assume a single market player, operating a 30 kWh battery with charging/discharging limits of 30 kWh without internal load or generation capabilities. Under static prices ( $\lambda^{buy} = \lambda^{feedin} = \text{Baseline price}$ ,  $\gamma = \beta = 0$ ), the optimal solution for the player is to charge fully at  $t = 2$  and discharge fully at  $t = 5$ , yielding a profit of 4.5\$. Given this behavior, the GSO pays a cost of 138.75\$.

Timestamp	Baseline Price (\$)	Base Demand (kWh)
1	0.40	40
2	0.35	35
3	0.40	40
4	0.45	45
5	0.50	60
6	0.45	45

Table 2: Baseline demand and prices

**Impact of Price Intervention** Now assume the ISO is willing to implement the intervention, and to set non-linear price signals. The ISO optimizes the price parameters, setting  $\beta = 0.002$  and  $\gamma = 0.455$  by solving for the objective function described in Equation ISO objective. This price adjustment incentivizes the player to redistribute charging and discharging activities, as the player solves the problem described in Section B.3. The optimal strategy for the player is as shown in Table 3, resulting in a higher profit of 6.52\$, including a subsidy from the ISO to the player (via sublinear feed-in price component) of 3.27\$. For the ISO total costs are reduced to 118.21\$ with the subsidy included. The intervention eliminates inefficiencies, benefiting both the ISO and the market player.

This example highlights the potential of price intervention to align market players' behavior with system-level efficiency goals. Furthermore, it demonstrates that the price intervention is not a zero-sum game, and some interventions can be beneficial for all parties involved.



Figure 1: Original demand

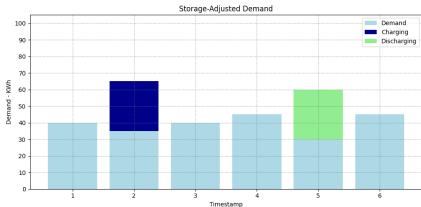


Figure 2: Linear Prices

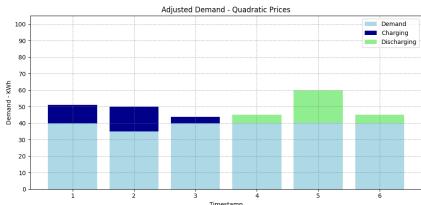


Figure 3: Quadratic Charging, Sublinear Discharging

Table 3: Non-linear prices implementation

However, what is described here is just one price intervention type possible. In general, the ISO would explore the space of all possible price interventions to find the optimal one. We suggest searching in this space using RL methods.

## F Day-Ahead Pricing as a Bandit Problem

At time 0, the ISO fixes prices in advance for all  $t$ , and receives a reward after the 48-timestep episode ends. This makes the ISO decide about the prices once per episode, which matches the dispatch decision. This turns the problem into a (very complex) bandit problem.

The bandit problem for dispatch and pricing in an electricity market is defined by the tuple:

$$\mathcal{B} = \langle \mathcal{A}, \mathcal{R}, \mathbb{P}, T \rangle$$

where:

- $\mathcal{A} = \{(d, p) \mid d \in \mathcal{D}, p \in \mathcal{P}\}$  is the set of actions, where each action is a pair  $(d, p)$ :
  - $d \in \mathcal{D}$ : Dispatch decision representing the amount of power to produce or allocate at a given time.
  - $p \in \mathcal{P}$ : Price levels, including selling prices and feed-in tariffs offered to market participants.
- $\mathcal{R}(d, p)$  is the reward associated with selecting the action  $(d, p)$ . Here, the reward is defined as the negative cost incurred by applying  $(d, p)$ , such that:

$$\mathcal{R}(d, p) = -C(d, p),$$

where  $C(d, p)$  represents the total operational cost, including dispatch costs, market costs, and reserve activation costs.

- $\mathbb{P}(d, p)$  denotes the probability distribution governing the outcomes (e.g., market responses, demand realization) associated with the action  $(d, p)$ .
- $T$  is the time horizon, representing the total number of decision rounds.

At each time step  $t \in \{1, 2, \dots, T\}$ , the agent selects an action  $(d_t, p_t) \in \mathcal{A}$ , observes the resulting market dynamics and incurred cost  $C(d_t, p_t)$ , and receives a reward  $\mathcal{R}(d_t, p_t) = -C(d_t, p_t)$ .

The objective is to minimize the cumulative cost over the time horizon  $T$ , minimizing the cumulative regret  $R_T$ , defined as:

$$R_T = \sum_{t=1}^T C(d^*, p^*) - \sum_{t=1}^T \mathbb{E}[C(d_t, p_t)],$$

where  $(d^*, p^*)$  is the optimal dispatch and pricing policy that minimizes the expected cost:

$$(d^*, p^*) = \arg \min_{(d, p) \in \mathcal{A}} \mathbb{E}[C(d, p)].$$

This formulation addresses the trade-off between exploration (testing new dispatch and pricing strategies to learn their outcomes) and exploitation (applying strategies believed to minimize costs based on current knowledge).

## G The Energy Market as MARL

In modeling modern power systems using marl, it is essential to account for multiple interacting perspectives. These include the physical constraints of the grid (e.g., stability limits), agent-level decision processes under partial observability, and the heterogeneity of demand profiles encompassing both nominal and flexible demand. Effective models must also incorporate market and pricing signals that influence agent behavior, and the temporal-spatial scalability required for real-world deployment. While these considerations are crucial for realistically and robustly capturing decentralized control strategies in complex energy environments, they also pose significant challenges to preserving the underlying Markovian structure that traditional agent-based decision models rely on.

### G.1 Formal Model

Through the lens of RL, the ISO aims to learn an optimal policy that balances overall system efficiency with the mitigation of risk, such as insufficient power supply and grid instability. Simultaneously, market participants seek to maximize their individual utility in response to market signals, subject to their own operational constraints and preferences. We formally

model this decentralized setting as a Markov Game (see Section 2), involving two types of agents: the ISO, and the GEAgents.

An important characteristic of the setting we aim to model is that the state space, action space, and reward functions are relatively straightforward to define. The complexity of solving this setting arises from modeling the joint transition function: the next state of the system and its stability depend on the actions performed by all agents.

### Modeling the ISO

- **State Space  $\mathcal{S}$ :** Every time step  $t$ , typically representing a half-hour interval, the system state is associated with a vector  $s_t \in \mathcal{S}$  that specifies operational factors that may affect decision-making. For the ISO this includes the system-level demand forecast  $\hat{D}$  for the specified horizon, the system-level realized demand  $D_t$  for the current time step, supply capacities, storage states, etc. It may also include factors that indicate the stability state of the system, for example, whether the supply-demand balance is violated.
  - **Action Space  $\mathcal{A}$ :** The ISO actions include the dispatch directives  $\Delta_t$  that are given for each time step  $t$  and setting the sell prices  $\xi_t(\cdot)$  and buy prices  $\phi_t(\cdot)$  for each time step. In real-time the ISO also activates reserves and curtails power if needed, but assume these actions are dictated by the state and require no decision-making.
- Importantly, we support two types of pricing dynamics. In a day-ahead pricing regime, the ISO makes the prices public at  $t = 0$ . In an online pricing setting, the ISO can dynamically set prices in response to the market signal. We discuss several pricing mechanisms and their characteristic, including the benefits of applying quadratic pricing, in Section 5.
- **Reward Function  $\mathcal{R}$ :** The ISO's reward integrates the economic efficiency and a risk measure to account for potential adverse outcomes arising from strategic GEAgents such that:

$$\mathcal{R} = -(C^{\text{dispatch}} + C^{\text{online}}) \quad (\text{ISO objective})$$

### Modeling the GEAgents

- **State Space  $\mathcal{S}$ :** Each GEAgent is associated with a PCS-unit for which the state includes its local information (e.g., state-of-charge) as well as the price signal advertised by the ISO.
- **Action Space  $\mathcal{A}$ :** Modern GEAgents have significant decision-making autonomy, allowing them to choose how much energy to store, consume, or sell based on their local goals, capabilities, and constraints. In this work, we assume the GEAgent sees the current prices and its local state at the start of each iteration before deciding how to act. Also, both generation and consumption are non-controllable. Specifically, we only support generation via pv and consumption that is part of the non-flexible load of the PCS-unit. This means that generation and production are exogenous to the agent and are governed by a stochastic process, and the only decision variable is the charge and discharge actions, which may have stochastic effects.
- **Reward Function  $\mathcal{R}$ :** For each GEAgent  $i$ , the step-wise reward is the net cash flow obtained by trading with the grid:

$$\mathcal{R}_t^i = \phi_t(P_t^s, P_t^b) - \xi_t(P_t^s, P_t^b).$$

Maximising the cumulative sum of  $\mathcal{R}_t^i$  over the horizon is equivalent to the strategic objective stated in (Strategic Player Objective), but written here without the expectation or the explicit time-index summation.

**Joint Transition Function  $\mathcal{T}$ : Influence of Multiple Agents:** Unlike a single-agent MDP, the Markov Game framework allows each agent's choice (including how GEAgents respond to prices or storage opportunities) to influence the next state. As mentioned above, the difficulty of modeling the transition function is at the core of the challenge. In general, the transition function can be decoupled into the state variables that are covered by the physical dynamics of the system. For example, when a charge or discharge action is performed, the battery dynamics obey (??); if an attempted action would violate  $0 \leq \text{SoC} \leq B_{\max}$

the short-fall or spillage is automatically settled with the grid, and a penalty is incurred. Propagated effect of local decisions, e.g., those solved with power flow.

Perhaps the most challenging aspect stems from the strategic interactions of the agents. These strategic decisions create a coupled system where each agent’s payoff depends on the actions of others. In principle, the Markov Function  $T(s_{t+1} | s_t, a_t^{ISO}, a_t^{PCS-unit})$  must fold together physical power flows, stochastic demand, renewables, battery chemistry and market clearing. Writing a closed-form  $T$  that captures all these layers is hopeless. Instead, we created the Energy-Net simulator (Section 6) maintain the physics and book-keeping, and we *learn* directly from roll-outs. This side-steps the need for explicit modeling of the complex dynamics and allows extracting value functions and policies using deep neural networks, rather than from first principles.

**Episode** As is typical in the day-ahead market, at the beginning of each episode (timestep  $t = 0$ ) the ISO receives the predicated demand  $\hat{D}$  for the next 48 half-hour intervals. It also receives the production and reserve capacities of its controlled units, the prices of each generated unit, and other information that might be relevant (i.e., weather forecast, special events, etc.). If day-ahead pricing is applied, the ISO sets and advertises the  $\xi_t(\cdot)$  and feed-in tariff  $\phi_t(\cdot)$  for the whole episode.

At each subsequent timestamp ( $1 \leq t \leq 48$ ), the following sequence of events occurs:

- (a) The ISO observes the realized demand  $D_t$ .
- (b) If dynamic pricing is applied, the ISO sets the sell price  $\xi_t(\cdot)$  and feed-in tariff  $\phi_t(\cdot)$  for timestamp  $t$ .
- (c) The GEAgents strategically respond to the prices by buying or selling power to the grid.
- (d) If the net demand after accounting for the net power  $P_t^{\text{net}}$  exceeds the scheduled production ( $\Delta_t$ ), the ISO activates reserves (e.g., peaker plants) to cover the shortfall or curtails power to cover overloads.

This iterative process continues until the end of the planning horizon. Both agents seek a stationary (possibly stochastic) policy that maximizes their own long-term discounted accumulated reward.

## H The Energy-Net Simulator

In spite of a variety of simulators that currently exist [], there is no current framework that allows modeling the complex structure we want to account for and that is designed to work with off-the-shelf rl and marl methods. We therefore develop a novel simulator, Energy-Net<sup>4</sup>, that we will use to examine our proposed solutions. Energy-Net is a modular, discrete-time simulator of a hybrid electricity market. The environment we develop is flexible and adaptable, and can be used to accommodate different system configurations. At the core of the design of the software is a decoupling between the physical dynamics of the electrical system and the strategic agents. Energy-Net is built around a strict *physics-agent split*. A high-fidelity physical core advances loads, renewables, batteries, and reserves, while the ISO and GEAgent interact only through a Gym-style `step()` interface. This design (i) lets us plug in any off-the-shelf rl/marl algorithm without touching the power-system code, (ii) isolates market rules in a single controller module, and (iii) ensures that learned policies can affect the grid *only* via explicit levers—prices and dispatch tweaks—thus preserving physical realism while streamlining experimentation.

Building on the formal setting introduced in Section ??, Energy-Net instantiates the 24-hour day-ahead electricity market. A single simulation episode therefore comprises  $T$  uniform intervals of length  $\Delta t$  (in our experiments  $T=48$  and  $\Delta t=30\text{min}$ ), together covering one 24-hour operational horizon. At each step  $t \in \{1, \dots, T\}$  the environment reveals the current forecast and grid state to the agents, applies their actions, propagates the physical dynamics, and returns next-state observations and rewards through the standard Gym `step` interface.

---

<sup>4</sup>link to repo - removed to respect the blind review process

## H.1 Physical Layer

**Demand.** System demand at each step is modelled as

$$D_t = f_{\text{seasonal}}(t) + \varepsilon_t,$$

where  $f_{\text{seasonal}}(\cdot)$  captures the deterministic daily profile and  $\varepsilon_t \sim \mathcal{N}(0, \sigma^2)$  is zero-mean Gaussian noise with user-configurable standard deviation  $\sigma$ .

**GEOAgent.** Every PCS-unit hosts a single-block battery whose state of charge obeys

$$\sigma_{t+1} = \sigma_t + \eta_c [a_t]_+ \Delta t - \eta_d^{-1} [-a_t]_+ \Delta t,$$

subject to  $0 \leq \sigma_t \leq S_{\max}$  and  $|a_t| \leq P_{\max}$ . Here  $a_t$  is the charge ( $> 0$ ) / discharge ( $< 0$ ) power,  $\eta_c, \eta_d$  are efficiency factors, and  $P_{\max}$  the power limit.

Besides storage, each unit experiences *stochastic local load*  $l_t$  and PV generation  $g_t$ , drawn from configurable distributions. The net exchange with the grid is therefore

$$P_t^{\text{net}} = a_t + g_t - l_t.$$

**Reserve.** If  $\Delta_t + P_t^{\text{net}} < D_t$ , spinning reserve is activated and the simulator logs the penalty  $C_t^{\text{reserve}}(D_t - \Delta_t - P_t^{\text{net}})$ , whose functional form and coefficients are user-configurable.

## H.2 Market Layer

At each step  $t$  the ISO broadcasts a **buy tariff**  $\phi_t(\cdot)$  (applied to energy flowing *into* storage) and a **sell tariff**  $\xi_t(\cdot)$  (applied to energy flowing *out of* storage). Energy-Net supports two pricing regimes:

- a) **Online linear.** The operator chooses two bounded scalars  $\lambda_t^{\text{buy}}, \lambda_t^{\text{sell}}$  and sets

$$\phi_t(P) = \lambda_t^{\text{buy}}, \quad \xi_t(P) = \lambda_t^{\text{sell}}.$$

- b) **Quadratic (super-/sub-linear).** At the beginning of each episode ( $t = 0$ ) the operator fixes four coefficients  $\{\lambda^{\text{buy}}, \lambda^{\text{feedin}}, \beta, \gamma\}$ ; they remain unchanged for all subsequent steps. Power-dependent tariffs are then computed with exactly the same notation used in Section ??:

$$\xi_t = \lambda^{\text{buy}} P_t^b + \beta [P_t^b]^2, \tag{1}$$

$$\phi_t = \lambda^{\text{feedin}} P_t^s + \gamma \sqrt{P_t^s}. \tag{2}$$

Here  $\beta$  adds a *super-linear* surcharge to purchases, whereas  $\gamma$  grants a *sub-linear* bonus on injections. Optional real-time dispatch perturbations  $\delta_t$  can still be issued on top of these pre-committed price curves.

### H.2.1 Agent Interfaces

**ISO observations.** At each step  $t$  the operator receives  $(t, \hat{D}_t, \hat{P}_t^{\text{net}})$ , where the hat denotes a one-step-ahead forecast of the aggregated exchange of all PCS-units.

**PCS observations.** Every storage unit observes the tuple  $(t, \xi_t, \phi_t, \sigma_t)$ .

#### ISO actions.

- *Online linear.* Set the instant tariff pair  $(\xi_t, \phi_t)$  (+ optional dispatch tweak  $\delta_t$ ).
- *Quadratic (super-/sub-linear).* Commit the coefficient quadruple  $(\lambda^{\text{buy}}, \lambda^{\text{feedin}}, \beta, \gamma)$  that parameterises (??); these remain fixed for the whole episode.

**PCS action.** A single continuous decision  $a_t \in [-D_{\max}, C_{\max}]$  interpreted as charge [ $a_t > 0$ ] or discharge [ $a_t < 0$ ].

### H.2.2 Reward Structure

Per-step rewards follow the definitions already introduced in Section B.4.

### H.2.3 Multi-Agent Execution

Energy-Net wraps both agents in a *single* multi-agent environment that extends the GYM-NASIM interface [24]. `step(...)` consumes a dictionary of actions and returns observation, reward, and termination tuples keyed by agent identity. Internally, a unified EnergyNetController advances the simulation in the following sequential order:

1. **Price setting** — the ISO chooses tariffs (and, if enabled, dispatch).
2. **Battery control** — the PCS-unit responds with its charge or discharge command.
3. **Energy exchange** — supply, demand, and storage flows are balanced; any shortfall triggers spinning reserve.
4. **State update and reward** — physical states, SoC, and financial ledgers are updated, and rewards are computed for both agents.

This integrated design eliminates manual data transfer between separate environments and exposes consistent, step-level metrics for training and evaluation. Notably, additional assets — renewables, alternative storage chemistries, custom reward definitions — can be introduced by registering new modules that comply with the interfaces above; no modification of the core simulation loop is required.

## I Evaluation Setup

Table 4: Scenario matrix used throughout Section ???. Columns B–C describe the **ISO** policy elements; column F the **PCS**. “Learned (prior  $S_k$ )” means the dispatch network is frozen from scenario  $S_k$  while the remaining degrees of freedom are (re-)trained. Every “Learned” block uses TD3.

ID	ISO pricing	ISO dispatch	PCS behaviour
Baseline	N/A	Equal to predicted demand	N/A
ISO-Dispatch	N/A	Learned	N/A
ISO-L	Online linear	Learned (prior S2)	Deterministic / fixed
ISO-Q	Quadratic	Learned (prior S2)	Deterministic / fixed
Joint-Storage-L	Online linear	Learned (prior S3)	<b>Learned</b>
Joint-Storage-Q	Quadratic	Learned (prior S3)	<b>Learned</b>
Joint-PCS-L	Online linear	Learned (prior S4)	<b>Learned + intrinsic load/production</b>

### Global scenario parameters (all baselines).

- Demand pattern: sinusoidal

$$D_t = L_0 + A \cos\left(\frac{2\pi}{P}(kt + \phi)\right)$$

with base load  $L_0 = 150$  MWh, amplitude  $A = 50$  MWh, interval multiplier  $k = 8$ , phase shift  $\phi = 5$ , period divisor  $P = 24$ .

- Dispatch energy price: \$100 per MWh.
- Reserve energy price: \$300 per MWh.
- Forecast-error noise (prediction error):  $\sigma = 10$  MWh.

For each interval  $t$  we first sample the *realised* demand  $D_t$  from the sinusoidal profile above. The ISO observes only a noisy one-step-ahead prediction

$$\hat{D}_t = D_t + \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, \sigma^2).$$

Hence, each experiment measures both the forecast error and the operator’s reaction to it. Note that even in the *day-ahead* pricing scenarios—where the six tariff coefficients chosen at  $t = 0$  remain fixed throughout the episode—the instantaneous ISO reward  $r_t^{\text{ISO}}$  is still computed *at every step*. This preserves time-resolved feedback while respecting the regulatory commitment to day-ahead prices.

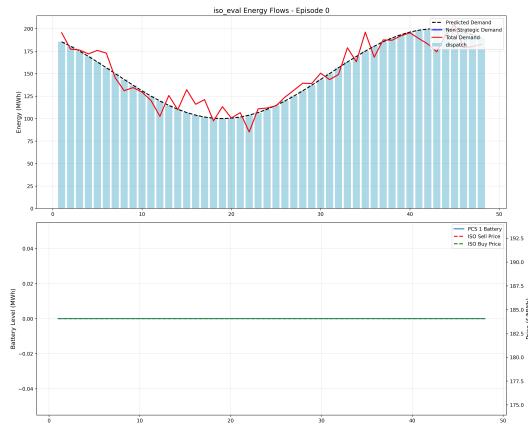


Figure 3:

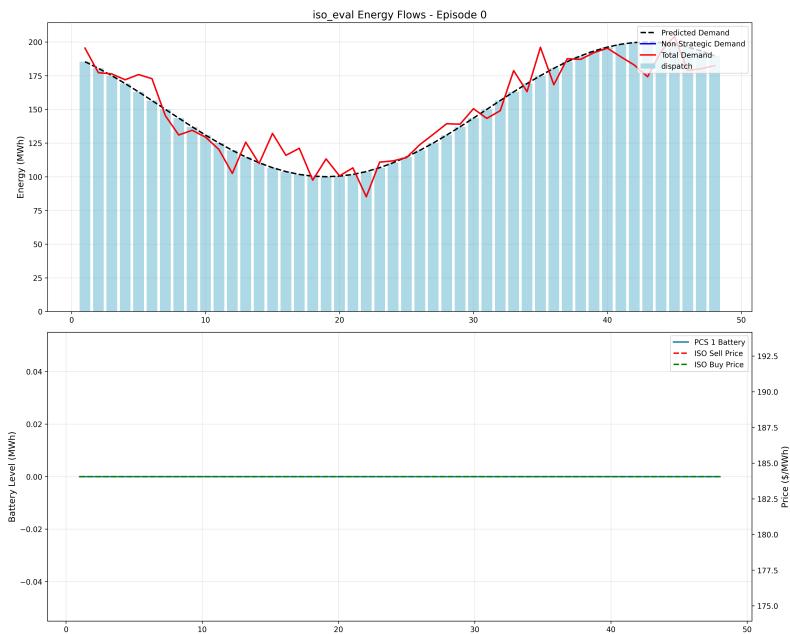


Figure 4: Energy, tariff, and cost traces for **S5-L**. Top: dispatch vs. realised demand. Middle: ISO buy/sell tariff trajectories. Bottom: cumulative cost distribution at episode end.

Table 5: Episode–total *cost* and *energy* breakdown across all evaluated scenarios (see Table 4 for scenario definitions).

Scenario	Dispatch		Reserve		PCS-unit Exchange
	Cost [\$]	Energy [MWh]	Cost [\$]	Energy [MWh]	Cost [\$]
ISO-L	728 235.04 ± 5089.24	7 282.34 ± 50.89	52 931	176.05 ± 19.07	10 276
ISO-Q	750 606	7 506	36 390	121	201
Joint-Storage-L	812 603	8 126	44 455	148	0
Joint-Storage-Q	812 600	8 126	44 400	148	0
Joint-PCS-L	732 244	7 322	50 466	168	12 801
Joint-PCS-Q	745 000	7 450	35 000	117	1 000

Table 6: Episode–total *cost* and *energy* breakdown across all evaluated scenarios (see Table 4 for scenario definitions).

2*Scenario	Dispatch		Reserve		PCS-unit Exch.
	Cost [\$]	Energy [MWh]	Cost [\$]	Energy [MWh]	Cost [\$]
ISO-L(Lin. fixed)	728 239	7 282	52 931	176	10 276
ISO-Q(Quad. fixed)	750 606	7 506	36 390	121	201
Joint-Storage-L(Lin. joint)	812 603	8 126	44 455	148	0
Joint-Storage-Q(Quad. joint)	812 600	8 126	44 400	148	0
Joint-PCS-L(Lin. joint+local)	732 244	7 322	50 466	168	12 801
Joint-PCS-Q(Quad. joint+local)	745 000	7 450	35 000	117	1 000
<b>Total</b>	—	—	—	—	—

## J Results

- (a) **Local context re-activates storage.** Without local load/generation (B4 scenarios) the TD3-trained ISO simply caps  $\xi_t$ , starving the battery of arbitrage opportunities. When local processes are introduced() the PCS–unit intervenes in roughly one-third of the time-steps, shaving peaks and flattening troughs. **Reserve capacity**
- (b) **Quadratic pricing yields the best balance.** The sensitivity-based quadratic tariff gives the lowest system-wide cost (\$781 k) with almost no cash transfer to the PCS-unit. Linear pricing still helps but pays the battery  $\approx \$13$  k to achieve similar savings.

### J.1 Empirical Evaluation Process

#### J.1.1 Baseline– Fixed Day-Ahead Schedule

Table ?? reports the episode-level costs when no real-time prices are published and the PCS-unit is inactive. These values serve as the *reference benchmark*; every subsequent result is expressed as a percentage change relative to the mean total cost reported here (\$782,408).

Table 7: Episode-total cost breakdown for Baseline

Component	Cost [\$]	Energy [MWh]	Share [%]
Dispatch	720,000	7,200	92.0
Reserve	62,408	208	8.0
PCS exchange	0	0	0.0
<b>Total cost</b>	<b>782,408</b>	<b>7,408</b>	100.0

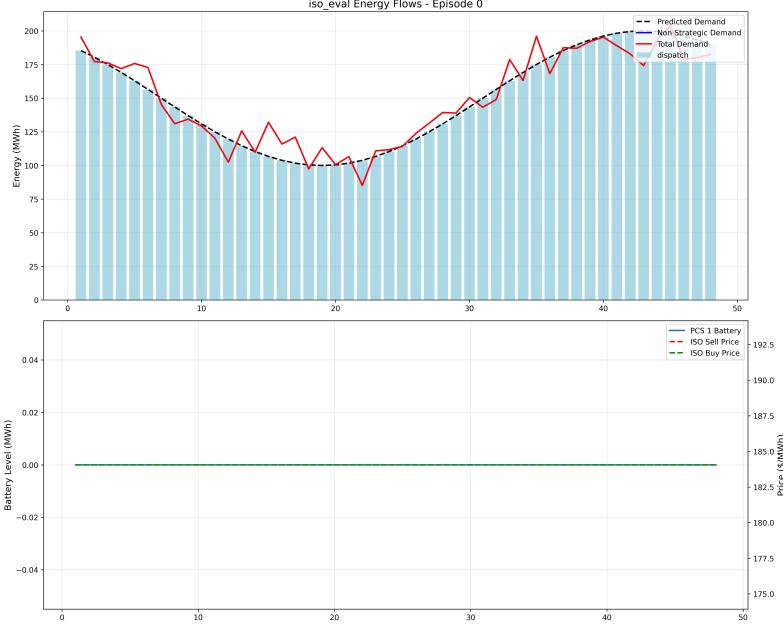


Figure 5: Energy–flow profile for Baseline. *Top*: dashed line = forecast demand, solid red = realised demand, blue bars = fixed day–ahead dispatch. *Bottom*: battery state of charge stays at 0 and buy/sell prices coincide at a constant level, confirming that no storage actions or dynamic tariffs are present.

### J.1.2 ISO-Dispatch – Adaptive Dispatch, No Price Signal

In this configuration the ISO is allowed to *revise the dispatch plan* every time step to follow its demand forecast, but still publishes no real-time tariffs; the PCS remains idle. Compared with Baseline, this setting isolates the value of feed-forward unit-commitment alone.

Table 8: Episode-total cost breakdown for ISO-Dispatch

Component	Cost [\$]	Energy [MWh]	Share [%]
Dispatch	722,881	7,229	90.6
Reserve	74,662	249	9.4
PCS exchange	0	0	0.0
<b>Total cost</b>	<b>797,543</b>	<b>7,478</b>	100.0

Relative to the fixed schedule (\SettingOne), total cost *increases* by +1.9 % because the higher dispatch level more than offsets the modest reserve reduction. This confirms that unit-commitment decisions alone do not capture real-time variability in a cost-effective way when no flexibility is available.

### J.1.3 ISO-L – Linear Price Signal with *Pre-defined* PCS Actions

The ISO now publishes real-time *online linear* buy and sell prices while continuing to adapt dispatch. The PCS does **not** react to those prices; instead it executes a fixed, offline–designed charging profile (charge in the early-valley hours, discharge during the evening peak), hence all storage actions are deterministic and price-agnostic.

Compared with the fully rigid benchmark (\SettingOne), reserve expenditure drops by –15 % ; the deterministic battery schedule absorbs part of the forecast error even without price responsiveness. The gain is almost entirely transferred to the PCS through the \$10k exchange payment, so total system cost changes by less than one percent.

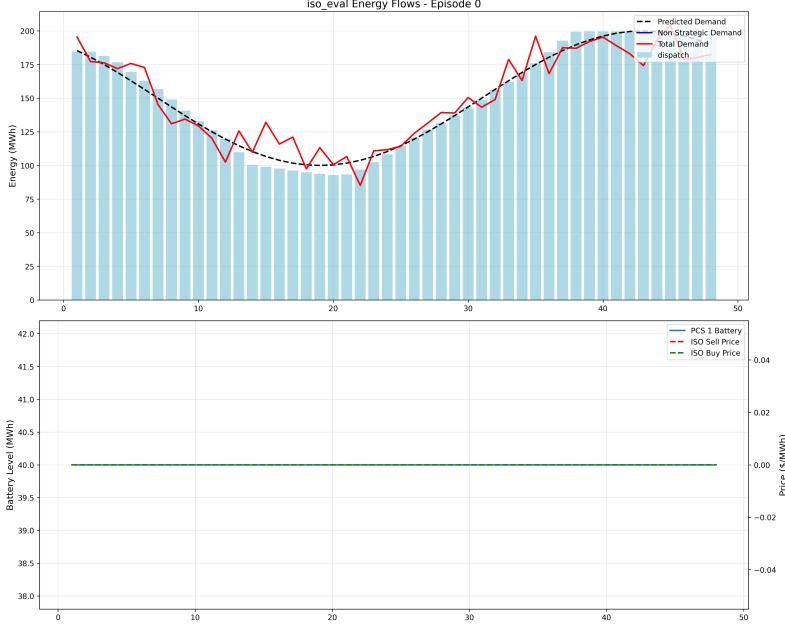


Figure 6: Energy-flow profile for ISO-Dispatch. *Top*: dashed line = forecast demand, solid red = realised demand, blue bars = fixed day-ahead dispatch. *Bottom*: battery state of charge stays at 0 and buy/sell prices coincide at a constant level, confirming that no storage actions or dynamic tariffs are present.

Table 9: Episode-total cost breakdown for ISO-L

Component	Cost [\$]	Energy [MWh]	Share [%]
Dispatch	728,239	7,282	92.0
Reserve	52,931	176	6.7
PCS exchange (paid)	10,276	800	1.3
<b>Total cost</b>	<b>791,447</b>	<b>7,458</b>	100.0

#### J.1.4 ISO-Q– Quadratic Price Signal with Pre-defined PCS Actions

Here the ISO adopts the *quadratic* pricing scheme (three coefficients for buy price and three for sell price) while keeping the same deterministic battery schedule used in ISO-L.

Table 10: Episode-total cost breakdown for ISO-Q

Component	Cost [\$]	Energy [MWh]	Share [%]
Dispatch	750,606	7,506	95.3
Reserve	36,390	121	4.6
PCS exchange (paid)	201	—	0.1
<b>Total cost</b>	<b>787,197</b>	<b>7,627</b>	100.0

**Effect relative to \SettingOne .** Quadratic pricing further suppresses reserve expenditure (−31% vs. \SettingOne and −31% vs. \SettingThreeL ) while almost eliminating ISO payments to the PCS (\$201 compared with \$10 k in \SettingThreeL ). The remaining cost is dominated by a higher dispatch level, so total system cost is −0.7% below the fixed benchmark and nearly identical to ISO-L.

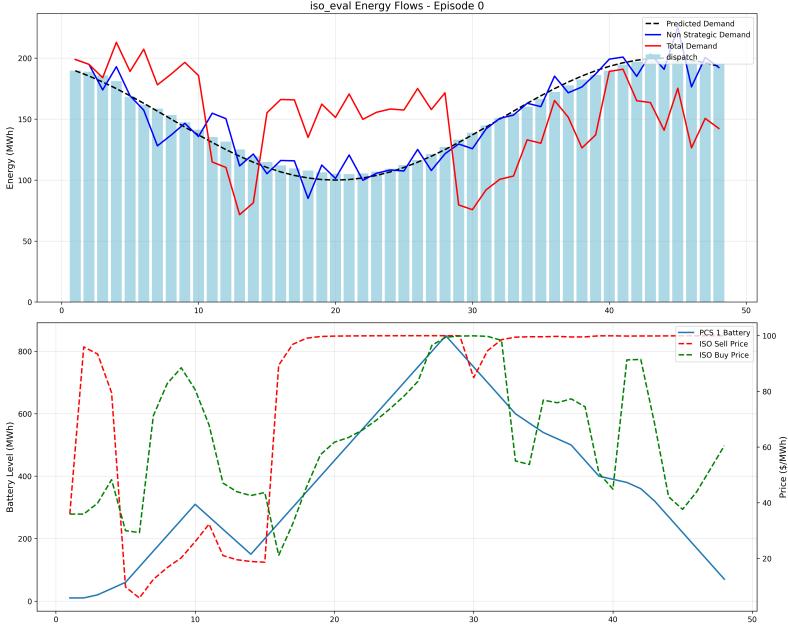


Figure 7: Energy–flow profile for ISO-L. A pre-scheduled battery cycle.

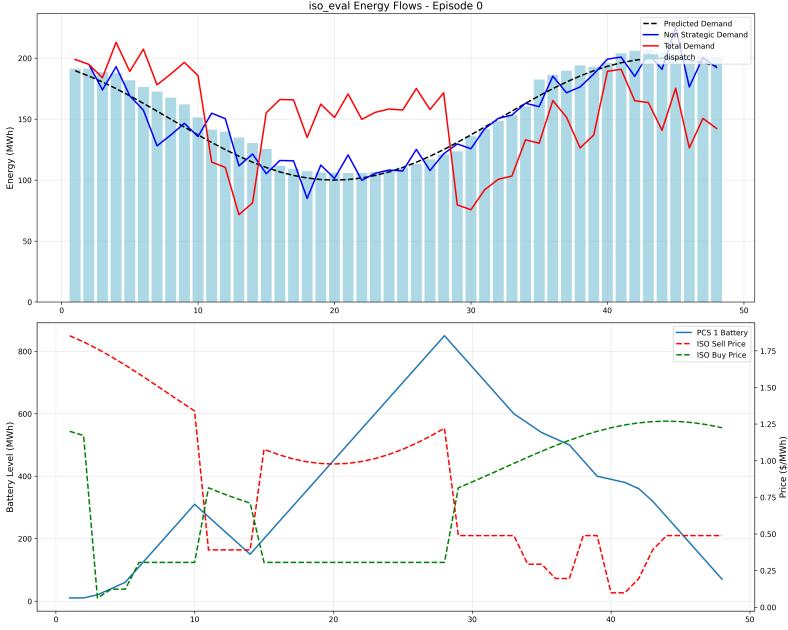


Figure 8: Energy-flow profile for ISO-Q(quadratic prices, deterministic PCS).

## J.2 Joint-Storage-L+ Joint-Storage-Q- TD3 ISO Monopolistic Pricing, Learned PCS

When both ISO and PCS are trained with TD3 under any pricing policy, the ISO rapidly discovers that it can *extract all surplus* by raising buy/sell tariffs to the upper bound, leaving no profitable arbitrage opportunity for the PCS. As a result, the battery remains inactive (zero net exchange), and the ISO effectively adopts a monopolistic pricing policy.

**Key observation.** The TD3-trained ISO converges to the maximum allowable tariffs, so the PCS agent—though fully capable—opts to remain idle, yielding zero net payments. This

Table 11: Episode-total cost breakdown for Joint-Storage-L.

Component	Cost [\$]	Energy [MWh]	Share [%]
Dispatch	812,603	8,126	94.9
Reserve	44,455	148	5.1
PCS exchange	0	0	0.0
<b>Total cost</b>	<b>857,058</b>	<b>8,274</b>	100.0

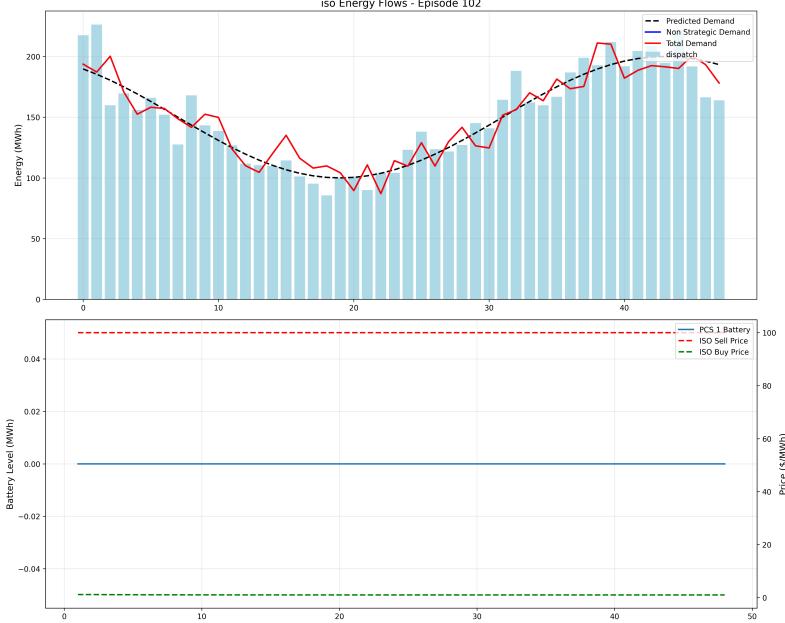


Figure 9: Joint-Storage-L. Under linear pricing, the ISO learns to set buy/sell tariffs at their upper bounds, extracting all potential profit from the PCS (zero net exchange) while achieving only modest reserve reduction.

“monopolistic” equilibrium eliminates storage activity entirely, demonstrating that without appropriate constraints or incentive alignment, pure RL pricing can lead to degenerate outcomes.

### J.3 Joint-PCS-L– Learned ISO and PCS under Endogenous Load & Production

In S4 the joint TD3 training converged to a *degenerate equilibrium*: once the ISO discovered that maximal tariffs eliminate any profitable arbitrage, the battery stayed idle and the market stalled. To re-introduce economic pressure we augmented the PCS with a small, uncontrollable background demand and a photovoltaic generation profile.<sup>5</sup> Whenever the combined load exceeded the instantaneous SoC the PCS had to *buy* from the grid, incurring a penalty; surplus solar could be fed back at the ISO’s sell price. This forces the storage unit to interact with the market even if arbitrage margins are slim.

---

<sup>5</sup>Both traces are purely illustrative: square-wave HVAC = 5kW peak; PV follows a bell curve peaking at 6kW around solar noon. Full details are provided in the supplementary material.

Table 12: Episode-total cost and energy breakdown for Joint-PCS-L.

<b>Component</b>	<b>Cost [\$]</b>	<b>Energy [MWh]</b>	<b>Share [%]</b>
Dispatch	972,866	9,729	99.3
Reserve	0	0	0.0
PCS exchange	<b>5,392</b>	341	0.5
<b>Total cost</b>	<b>978,928</b>	<b>9,729</b>	100.0