

Reinforcement Learning-Driven Local Transactive Energy Market for Distributed Energy Resources

Steven Zhang^a, Daniel May^a, Mustafa Gül^b, Petr Musilek^{a,c,*}

^a*Electrical and Computer Engineering, University of Alberta, Canada*

^b*Civil and Environmental Engineering, University of Alberta, Canada*

^c*Applied Cybernetics, University of Hradec Králové, Czech Republic*

Abstract

Local energy markets are emerging as a tool for coordinating generation, supply, storage, and consumption of energy from distributed resources. In combination with automation, they promise to provide an effective energy management framework that is fair and brings system-level savings. The cooperative-competitive nature of energy trading calls for automation based on multi-agent systems with learning energy-trading agents. However, depending on the dynamics of the agent-environment interaction, this approach may yield unintended behavior of market participants. Thus, the design of local energy market suitable for reinforcement learning agents must take into account this interplay. This article introduces autonomous local energy exchange as an experimental framework based on multi-agent approach combined with a double auction market mechanism. In this setting, participants determine their internal price signals and make energy management decisions through market interactions, rather than relying on predetermined external price signals. The market properties suitable for use with reinforcement learning are determined through simulation experiments. The results show that market truthfulness is required to maintain demand-response functionality, while weak budget balancing is needed to provide a strong reinforcement signal for the learning agents. The ability of the market to facilitate demand response is illustrated by comparing the resulting behavior with two baselines: net billing and time-of-use rates. The market-based pricing was significantly more responsive to fluctuations in the community net load. In addition, for the test community, the more accurate accounting of renewable energy usage reduced bills by a median 38.8% compared to net billing.

Keywords: Transactive energy, Demand response, Distributed energy resources (DER), DER integration, Local energy market, Reinforcement learning

1. Introduction

1 Demand response (DR) techniques have become popular means to increase the
2 value of distributed energy resources (DER), such as rooftop solar, while mitigating
3

*Corresponding author

4 the negative effects of their intermittent nature. DR methods can be direct or indirect.
5 Indirect DR aims to change customer behavior using an incentive signal, usually
6 through monetary means [1, 2, 3]. Direct DR grants grid operators immediate control
7 to perform grid balancing. As DER adoption continues, centralized approaches to DR
8 will encounter scalability barriers [4, 1, 5, 6], and more granular and robust control
9 will be necessary due to the increased intermittency and stochasticity of supply and
10 demand. Despite the effort to tackle some of these challenges [7, 8, 9], it is increasingly
11 clear that alternative, decentralized solutions, must be explored. In this context,
12 transactive energy (TE) is gaining popularity as a design framework for decentralized
13 DR [5, 6]. The U.S. Department of Energy Gridwise Architecture Council defined TE
14 as “a system of economic and control mechanisms that allows the dynamic balance
15 of supply and demand across the entire electrical infrastructure using value as a key
16 operational parameter” [10]. A locally constrained TE system is often referred to
17 as local energy market (LEM). In lieu of a formal definition of LEM [11, 12], we
18 adopt the description developed by Mengelkamp et al. [11], i.e. “a market platform
19 for trading locally generated (renewable) energy among residential customers within
20 a geographically and socially close community. Supply security is ensured through
21 connections to a superimposed energy system.”

22 High DR participation rates, especially within a LEM, can only be maintained
23 with automation. Expert-designed, rule based systems have initially been consid-
24 ered for this purpose. Learning-based approaches such are now preferred, mainly
25 because of their robustness and scalability. However, the vast majority of existing
26 approaches that apply learning methods to LEM do not tailor the market mecha-
27 nism to the algorithm used for automation. This is especially problematic given the
28 fact that most established LEM mechanisms were designed for human participants
29 or rule-based system automation [13]. This results in suboptimal DER utilization,
30 as the LEM is not appropriately adjusted to best leverage the potential of the au-
31 tomation approach used. This is exacerbated for reinforcement learning (RL) agents
32 that can quickly learn to exploit loopholes in competitive-collaborative multi-agent
33 settings [14]. Mengelkamp et al. clearly identify this research gap in their review,
34 stating that “a comprehensive comparison of the impacts of different trading designs
35 (especially market mechanisms) should be carried out. Specifically, the impact of
36 different allocation mechanisms on the market objectives and agent behavior need to
37 be evaluated” [11]. Their later work [15] follows the same reasoning, noting that both
38 agent design and market design influence the resulting system behavior. We argue,
39 that the LEM mechanism should be tailored to adequately empower the strengths of
40 their automation methods and to mitigate their potential weaknesses.

41 To the best of our knowledge, there has been no contribution that explicitly designs
42 experiments to identify the requirements that a specific LEM market mechanism must
43 fulfill to be well compatible with RL actors. This article aims to provide a starting
44 point to fill this research gap. We narrowly focus on the following two questions:

- 45 • What are the required properties of LEM settlement mechanism suitable for
46 deployment of RL-based automation?
- 47 • Does the resulting market behavior effectively support DR for LEM with high

48

penetration of DER?

49

To answer these questions, three different settlement mechanisms are examined that cover a set of established criteria for auction environments. The most suitable market design is found by analyzing the agent policies developed for each mechanism. This study then models the resulting LEM transactions as a dynamic price signal and compares its economic performance with existing pricing methods.

54

55

56

57

58

59

60

61

62

63

This article is organized in five sections. Section 2 provides the necessary background and describes the related work. Section 3 introduces the proposed autonomous local energy exchange (ALEX) and describes it as a stochastic game. Section 4 describes two sets of experiments. The first set is designed to identify settlement mechanism suitable for market automation using learning agents. The second set performs an economic analysis of the selected mechanism and compares its performance with several benchmarks. Major conclusions are summarized in Section 5, along with possible directions for future work. Three appendices describe the principles of net billing (Appendix A), transactive energy simulator T-REX (Appendix B), and details of the specific market design used in this article (Appendix C).

64

2. Background and Related Work

65

66

67

68

69

This section provides a brief introduction to RL and more in-depth review of the related work, focusing on articles that combine RL with LEM. For a broader context of LEMs, the reader may refer to a general review by Mengelkamp et al. [11], game-theory focused review by Pilz et al. [12], and review of LEM settlement and market mechanisms by Khorasany et al. [13].

70

2.1. Reinforcement Learning

71

72

73

74

75

76

77

78

79

There have been numerous approaches used to address the optimization problems inherent to demand side management [1, 4, 5, 6], including mixed-integer programming, stochastic programming, and dynamic programming. After RL demonstrated great competence in partially observable, stochastic game environments [16, 17, 18], this model-free, sequence-oriented, semi-supervised machine learning framework has also gained popularity as a control method for DR [19]. The learned policy can substitute the solution of the equivalent optimization problem at each time step. As a result, RL approaches are more computationally efficient at scale, when compared to conventional optimization methods.

80

81

82

83

In the RL setting, illustrated in Fig. 1, an agent learns to maximize the return G by interacting with its environment through actions a while receiving observations of the environmental state s . G is commonly defined as the expected, discounted cumulant of future reward R

$$G_t = \sum_{i=t}^T \gamma^{(i-t)} R_i, \quad \forall \gamma \in [0..1], \quad (1)$$

84 where γ is the discount factor. The expected value of G_t , given state s_t or the state-
 85 action tuple (s_t, a_t) , is referred to as the state value

$$V(s_t) = \mathbb{E}(G_t | s_t, \pi(s_t)), \quad (2)$$

86 or action value

$$Q(s_t, a_t) = \mathbb{E}(G_t | s_t, a_t), \quad (3)$$

87 respectively. An RL agent acts according to a policy

$$\pi : A \times S \rightarrow [0...1]., \quad (4)$$

88 Policy is a (probabilistic) mapping of the state space S on action space A

$$\sum_a \pi(a, s_t) = 1. \quad (5)$$

89 This allows the definition of the state value V as action value Q weighted by π

$$V(s_t) = \frac{1}{n_a} \sum_a \pi(a, s_t) Q(s_t, a). \quad (6)$$

90 This system of equations (1-6) is sufficient to broadly classify all RL algorithms
 91 along two axes: the learned function and the relation between the target and behavior
 92 policy. According to the learned function, RL algorithms can be classified as policy-
 93 gradient methods and value-based methods. The policy-gradient methods directly
 94 learn the policy π , while the value-based methods learn estimations for either V
 95 or Q , and employ a fixed mapping of these values to π . RL algorithms can also
 96 be classified into on-policy and off-policy methods, by comparing their target and
 97 exploratory behavior. An on-policy RL algorithm explores the environment with the
 98 same policy that is optimized, while an off-policy algorithm explores the environment
 99 with a behavioral policy $b \neq \pi$.

100 Internally, RL algorithms often employ function approximation techniques to per-



Figure 1: Reinforcement Learning Setting

101 form the mapping of the state space S to the learned target, and therefore the return
 102 G . As a framework, RL is independent of the choice of state estimator. Currently, a
 103 very popular choice is the use of deep artificial neural networks. Historically, other
 104 function approximation techniques have also been used, such as tabular encoding.

105 The algorithm used in this paper is Q -learning, a well-established, value-based,
 106 off-policy algorithm. It learns the greedy policy, a deterministic policy that always
 107 picks a corresponding to the largest Q , by following an arbitrary annealed behavioral
 108 policy b . A popular choice for b is the ϵ -greedy policy, which takes a random action
 109 with probability ϵ and otherwise follows the greedy policy. The corresponding learning
 110 rule can be written as

$$Q^{\text{updated}}(a_t, s_t) \leftarrow (1 - \alpha)Q(a_t, s_t) + \alpha \left(R_t + \gamma \max_a (Q(a, s_{t+1})) \right), \quad (7)$$

111 where α is the learning rate.

112 Q -learning is a relatively well-understood RL algorithm, with strong convergence
 113 criteria for the tabular function approximation case, as long as both ϵ and α are
 114 annealed towards 0 at infinity. It is also the most common algorithm in the related
 115 literature reviewed in section 2.2.

116 2.2. Reinforcement Learning for Local Energy Markets

117 Several authors investigate the combination of RL and dynamic pricing for central-
 118 ized control. Notably, Kim et al. [2], and Lu et al. [3] develop RL-based approaches
 119 for dynamic pricing from the perspective of a service retailer. Both articles address
 120 difficulties of predicting participant response to a pricing schedule by mitigating the
 121 reliance on accurate customer side information. A Markov decision process is formu-
 122 lated based on customer behaviour models and preferences. A Q -learning agent is
 123 trained to simultaneously minimize customer costs and maximize the service provider
 124 benefit. The two approaches differ in the formulation of the reward function, which
 125 is a major influencing factor in RL algorithms, in general. Lu et al. [3] use a weighted
 126 sum of retailer and customers, while Kim et al. [2] use a modelled utility function. Al-
 127 though both proposed approaches successfully implement dynamic pricing strategies
 128 without scheduling, they still rely on modeling consumer behavior and preferences
 129 via utility functions.

130 Zhang et al. [20] train an RL agent to manage a community-shared battery and
 131 trade its resources on a TE market to maximize economy. Thus, the reward function
 132 is the economic performance of the battery. The authors show that positive economic
 133 benefits can be achieved, even when considering the running costs of the battery.

134 Xiao et al. [21] investigate optimized trading between a large number of intercon-
 135 nected microgrids using a deep Q -network (DQN) based RL agent. Similarly to Kim
 136 et al. [2], a utility function is used to determine the reward. As expected, the more
 137 sophisticated DQN algorithm outperforms the benchmark hotbooting Q -learning al-
 138 gorithm.

139 Foruzan et al. [22] investigate the behavior of self-interested Q -learning agents,
140 exchanging energy within a microgrid through a LEM. The agent’s goal is to maximize
141 its own profit. Managed DERs include battery energy storage systems, rooftop solar,
142 wind and diesel generators. The participants’ stochastic behavior is approximated
143 using random models. The authors investigate several micro grid configurations and
144 perform an in-depth hyperparameter study of the RL algorithm with respect to return,
145 self-sufficiency and fairness.

146 Zhou et al. [23] combine a fuzzy rule-based system with Q -learning to train agents
147 to exchange energy resources over a peer-to-peer LEM setup whose pricing is directly
148 tied to the ratio of supply and demand. The authors investigate the performance of
149 several community configurations with ranging number of battery energy storage sys-
150 tems and renewable generation assets. They show that such a system setup generally
151 achieves lower bills than TOU and net-billing baselines.

152 Chen et al. [24] employ a DQN variant to automate the interactions of prosumers
153 equipped with battery energy storage system in a LEM. The RL agents’ action space
154 consists of four distinct, discrete actions covering buy/sell and charge/discharge oper-
155 ations. The learned policy surpasses an intuitive, rule-based strategy. It also outper-
156 forms a pure random policy equivalent to a zero-intelligence agent, originally proposed
157 by Ghode et al. [25] as a baseline for agent competence in automated markets. In
158 another article, Chen et al. [26] investigate the function of Q -learning based energy
159 brokers as LEM consensus mechanism for settlements with profit used as the agent’s
160 reward. Using several ablation and sensitivity studies, the authors show that the
161 brokers efficiently learn how to maximize their own profit and the efficiency of the
162 market.

163 Kim et al. [27] extend a DQN variant designed for stock trading applications
164 to manage participation of a household in peer-to-peer energy exchange. Using ex-
165 periments with different rate schemes, they demonstrate that the developed agent
166 outperforms the simplified versions in terms of loss minimization and revenue maxi-
167 mization.

168 Bose et al. [28] focus on emerging participant interaction within a fixed LEM setup
169 under differing levels of DER penetration. The authors demonstrate that RL-based
170 agents in such environment can cause partial energy self-sufficiency to emerge. They
171 also show that the degree of self-sufficiency and the complexity of agent interactions
172 depends on the level of DER penetration within the market.

173 Mengelkamp et al. [29] study three different extensions of the Erev-Roth RL algo-
174 rithm applied to automate LEM participation. They find that the extensions further
175 increase the self-sufficiency of the LEM when compared to the original Erev-Roth
176 algorithm [30].

177 Mengelkamp et al. [15] compare a peer-to-peer LEM against a closed book, double-
178 auction LEM with settlement rounds. The authors compare the performance of zero-
179 intelligence agents and “intelligent” agents adopted from Nicolaisen et al. [31] on both
180 LEM designs. They show that all market scenarios offer similar economic advantages,
181 with the peer-to-peer LEM used by intelligent agents slightly outperforming the other
182 variants. However, they also note that using one strategy on different markets results
183 in different price trends, The authors eventually conclude that agent strategy and

184 market design need to be co-developed to guarantee the system’s performance.

185 3. ALEX: Autonomous Local Energy Exchange

186 As shown in Section 2, the need to model customer behaviour via utility functions
187 and heavy reliance on forecasts may hinder the robustness and scalability of tradi-
188 tional DR techniques. Furthermore, due to the amount of DERs that are expected to
189 be on-line in the near future, certain infrastructure requirements may pose additional
190 barriers for the effective deployment of DR systems. For example, the communication
191 and computational costs for centralized control may grow too expensive, especially
192 for high temporal resolutions necessary to fully capture the intermittent behaviour of
193 rooftop solar panels and stochastic use patterns of electric vehicles.

194 To avoid these difficulties, this article proposes a distributed, multi-agent approach
195 combined with a double auction market mechanism. When designing this approach,
196 dubbed ALEX (autonomous local energy exchange), the following assumptions have
197 been made:

- 198 • Participants are self-interested and, therefore, prioritize their own economic
199 well-being in the decision making process.
- 200 • Participants are willing to defer some decision making regarding interactions
201 with indirect DR measures to automation (e.g., using RL agents).
- 202 • Each participating unit is equipped with a smart meter, and a sufficient amount
203 of high-resolution historical data is available to train the RL agents.
- 204 • The large-scale electricity grid that customers are connected to is an infinite
205 bus.

206 3.1. Core Concept

207 Conceptually, ALEX is a behind-the-meter DR technique for a localized commu-
208 nity using a double auction market as a coordination mechanism. Market participants
209 are the customers who live within the community. However, this could be expanded
210 to include entities that are only temporarily present, such as electric vehicles. The
211 market employs double auctions with a fixed settlement frequency Δt . For each in-
212 terval $[t, t + \Delta t)$, participants can communicate their intention to trade energy by
213 submitting bids

$$\text{bid}_t = (q_t^{\text{bid}}, p_t^{\text{bid}}), q \in [0 \dots q_{\text{max}}^{\text{ask}}], p \in [p_{\text{min}}, \dots p_{\text{max}}], \quad (8)$$

214 and asks

$$\text{ask}_t = (q_t^{\text{ask}}, p_t^{\text{ask}}), q \in [0 \dots q_{\text{max}}^{\text{ask}}], p \in [p_{\text{min}}, \dots p_{\text{max}}], \quad (9)$$

215 where bid_t and ask_t communicate the intention to buy or sell energy, respectively.
216 They are represented by tuples consisting of the desired quantity q and desired price
217 p of energy to be exchanged. Quantities are expressed in Wh, and can range from 0
218 to a designated maximum. $q_{\text{max}}^{\text{ask}}$ should be set to accommodate expected maximum

219 generation derived from historical data. Likewise, q_{\max}^{bid} should be set to accommodate
 220 the expected maximum load demand. Similarly, prices in \$ are within a designated
 221 window between p_{\min} and p_{\max} . Bids and asks are settled pairwise at the end of each
 222 settlement round, returning a settlement signal m_t to each participant

$$m_t = (q_t^{\text{settlement}}, p_t^{\text{settlement}}) \quad (10)$$

223 More details on the market implementation are provided in Appendix C. It is
 224 important to stress that in this setting, participants both determine the price signals
 225 and make energy management decisions through market interactions, whereas most
 226 other DR approaches simply have agents react to external price signals.

227 The settlement signal m_t is represented as a list of tuples containing the settled
 228 quantities, and the respective prices. It is important to note that participants only
 229 receive information about their settlements and, therefore, do not have access to
 230 information on the behaviour of other participants.

231 After the internal trades are concluded, any excess generation/demand within the
 232 community is exchanged with the electricity grid at retail prices. In this article, we
 233 assume net billing, a commonly used practice where excess energy is sold to the grid at
 234 price $p_{\text{sell}}^{\text{grid}}$ and deficient energy is purchased from the grid for price $p_{\text{buy}}^{\text{grid}}$ that includes
 235 fees. The behind-the-meter setup grants the community a window of profitability by
 236 deferring fees. This naturally bounds the range of internal market prices as follows

$$p_{\min} = p_{\text{sell}}^{\text{grid}} \leq p_{\text{market}} \leq p_{\text{buy}}^{\text{grid}} = p_{\max}, \quad (11)$$

237 If the interactions between participants are dominated by the law of supply and
 238 demand, then ALEX functions as a decentralized, indirect DR tool. Its pricing sched-
 239 ule is strongly correlated with the ratio of supply and demand within the market. This
 240 naturally provides economic incentives for all market participants to balance supply
 241 and demand. We demonstrate this using the experiments described in Section 4.1.

242 3.2. ALEX as a Stochastic Game

243 To analyze the properties of the proposed approach, the auction and strategic
 244 bidding by the actors can be described using a suitable mathematical model. A
 245 game-theoretic representation of ALEX can be derived by modelling the interactions
 246 of participants as a discounted stochastic game

$$\Gamma := (n, L, S, A, P, R) \quad \forall t \in [0..T], \lambda \in (0..1), \quad (12)$$

247 where n is the number of players, L is the list players of length $|L| = n$, S is the state
 248 space, A is the action space, P represents the state transition probabilities, R is the
 249 reward function, t is the current time step over the modelling period $[0..T]$, and λ is
 250 the discount factor.

251 Both S and A can be decomposed into n individual components S^i and A^i , as
 252 shown below

$$S = S^1 \times \dots \times S^n, \quad (13)$$

253

$$A = A^1 \times \dots \times A^n. \quad (14)$$

254 Superscript i refers to a specific individual L^i , while the subscript is reserved for
 255 time t . Note that the action space A is separated in notation from a specific set of
 256 actions a_t at time step t , as in the commonly used RL nomenclature introduced in
 257 Section 2.

258 State transition probabilities are defined for any set of actions a_t taken at time
 259 step t , as follows

$$\forall a_t : P(S_{t+1}|S_t, a_t) := S_t \rightarrow S_{t+1} \quad (15)$$

260 Analogous to the RL setting, the reward or payoff in the stochastic game at time
 261 step t is defined by

$$R_t := S \times A \rightarrow r, \quad (16)$$

262 which maps from (S_t, a_t) to a real number $r \in \mathbb{R}$. Similarly, each agent aims to max-
 263 imize their own return G_t (1). Thus, all participants use their individually developed
 264 policy π^i (4 and 5), to determine action set a_t^i based on observations from S_t^i .

265 At each time step, all agents can interact with the market by submitting bids (8)
 266 and asks (9). This leads to the following definition of action

$$a_t^i = (\text{bid}_t^i, \text{ask}_t^i, e_t^i), \quad (17)$$

267 where the additional parameter, e_t^i , is reserved for future expansion of the model, e.g.,
 268 to define nonmarket actions, such as battery management or thermal load control.

269 Finally, the state observations for each agent are defined as follows

$$S^i = (d_t^i, g_t^i, m_{t-1}^i), \quad (18)$$

270 where d_t^i and g_t^i are, respectively, the load demand and generation at time t , and m_{t-1}^i
 271 are settlements received at time $t - 1$.

272 Note that the transition probabilities P result from the collective actions of all
 273 agents. However, due to the pairwise settlement mechanism, market design, and the
 274 observation space, P is not fully accessible to L^i . This ensures that the developed
 275 model is a truly stochastic game.

276 At least one stable Nash equilibrium is guaranteed to exist within Γ , as long as n , S
 277 and A are finite. This condition can be guaranteed by limiting prices p to a reasonable
 278 decimal place accuracy (e.g., 4 or 5 significant digits commonly used in banking). A
 279 is logically bounded by the condition previously defined by (11). As a result, S must
 280 also be finite, and therefore each implementation of ALEX is guaranteed to exhibit
 281 at least one stable Nash equilibrium.

282 3.3. Automation using Reinforcement Learning

283 Since the interaction through the developed stochastic game Γ requires strategic
 284 competence, automating the interactions of participants (typically prosumers, but
 285 theoretically any grid-connected entity) with the LEM is a reasonable response to
 286 the difficulties of accurately modeling customer behaviour. The proposed approach

287 centers around training RL agents to perform market interaction and energy man-
 288 agement actions, compensating for nonoptimal human behaviour. RL is theoretically
 289 very suitable for this task, as the stochastic game described in the previous subsec-
 290 tion is equivalent to a Markov decision process under an established set of criteria,
 291 outlined in [32]. The framework developed in this section is set up to be algorithm
 292 agnostic. However, the subsequent experiments described in Section 4.1.1 employ
 293 independent Q-learning.

In this article, the reward function, R^i , is formulated as follows

$$R_t^i = \text{cost}_t^{i,\text{LEM}} - \text{profit}_t^{i,\text{LEM}} + \text{cost}_t^{i,\text{grid}} - \text{profit}_t^{i,\text{grid}} + c, \quad (19)$$

where,

$$\text{cost}_t^{i,\text{LEM}} = q_t^{i,\text{settled-bids}} \times p_t^{i,\text{settled-bids}}, \quad (20)$$

$$\text{profit}_t^{i,\text{LEM}} = q_t^{i,\text{settled-asks}} \times p_t^{i,\text{settled-asks}}, \quad (21)$$

$$\text{cost}_t^{i,\text{grid}} = q_t^{i,\text{grid-buy}} \times p_t^{i,\text{grid-buy}}, \quad (22)$$

$$\text{profit}_t^{i,\text{grid}} = q_t^{i,\text{grid-sell}} \times p_t^{i,\text{settled-sell}}, \quad (23)$$

and,

$$q_t^{i,\text{settled-bids}} + q_t^{i,\text{grid-buy}} = d_t^i \text{ (load demand)}, \quad (24)$$

$$q_t^{i,\text{settled-asks}} + q_t^{i,\text{grid-sell}} = g_t^i \text{ (generation)}. \quad (25)$$

294 Other considerations, such as social welfare costs, are explicitly excluded as the
 295 current aim is to study participant and system behaviour using pure economic perfor-
 296 mance. Nevertheless, these factors may be included in future studies. Since RL agents
 297 are trained using high-frequency smart meter data, explicit customer behaviour mod-
 298 els can be omitted, as a sufficient amount of data can better capture nuanced and
 299 individualized customer behavioural patterns.

300 4. Experiments and Discussion

301 4.1. Suitability of Settlement Mechanism

302 Fig. 2 illustrates a typical participant, represented by a prosumer's home. The
 303 home may contain any combination of generation, storage, and controllable loads. For
 304 the experiments, the generation and load sources are taken directly from smart meter
 305 data. As the focus of this article is to study the LEM's behavior under differing ratios
 306 of available supply and demand, profile shaping via load shifting or battery storage
 307 are not considered.

308 4.1.1. Experimental Design

309 This experiment focuses on the first research question: What are the required
 310 properties of LEM settlement mechanism suitable for the deployment of RL-based

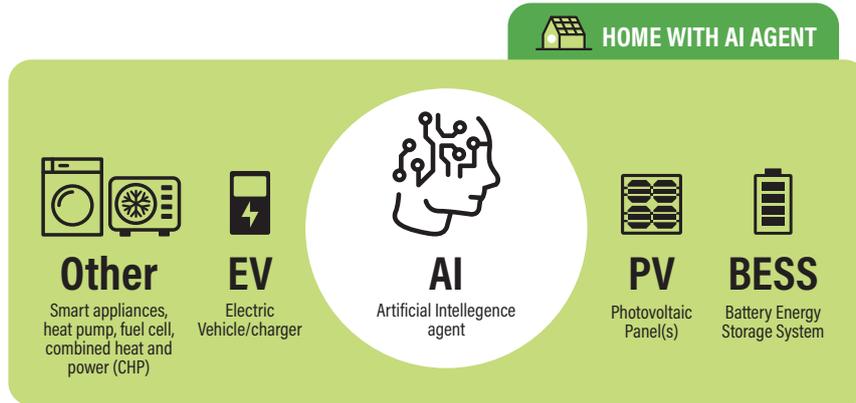


Figure 2: Diagram showing a typical participant, represented as a prosumer’s home

311 automation? The results of this experiment will also show how different market
 312 properties influence the policies learned by the RL agents. Since the environment
 313 (i.e., the market) plays just as important role as the learning algorithm, it is imper-
 314 ative to establish the most suitable market mechanism for subsequent research and
 315 implementation of more complex agent designs, and action and strategy spaces. The
 316 scope of the experiment is carefully managed to magnify the influence of the settle-
 317 ment mechanism on the resulting agent policies while providing strong convergence
 318 bounds despite ALEX’s properties as a partially observable, nonstationary environ-
 319 ment. As a reminder, the market design used in this study and the rules of interaction
 320 are described in detail in Appendix C. In this experiment, three different settlement
 321 mechanisms with varying market properties are tested:

- 322 1. Average-Price (M1): Trades are settled if the bid price is greater than or equal
 323 to the ask price. The settlement price is the average of the bid and ask prices.
- 324 2. Exact-Match (M2): Sellers and buyers choose bid and ask prices from a list of
 325 available prices. Trades are settled if the bid price equals the ask price.
- 326 3. Exact-Price (M3): Trades are settled if the bid price is greater than or equal
 327 to the ask price. The buyer buys from the auctioneer at the bid price, and the
 328 seller sells to the auctioneer at the ask price.

329 Any double auction mechanism can be described by the following properties [33]:
 330 individual rationality¹, budget balancing², truthfulness³, and economic efficiency⁴.

¹Individual rationality states that no participant should lose money from joining the auction

²There are two variants of budget balancing: weak and strong. In a weak budget balancing system, a portion of the money transferred also goes to the auctioneer; this is in addition to money transfers between participants which are the only type of exchange in a system with strong budget balancing.

³The dominant strategy in a truthful market is for the participants to report prices at what they believe should be the true value of the item to be exchanged.

⁴In an economically efficient system, at the end of all trading, the items should be in the hands of participants who bid the highest value.

331 An ideal mechanism satisfies all four properties, but it cannot be realized in practice.
 332 Since the design of ALEX and the use of RL agents ensures economic efficiency
 333 and individual rationality, the three settlement mechanisms can be differentiated by
 334 truthfulness and budget balancing alone, as shown in Table 1.

Mechanism	Market Property Settings			
	Individual rationality	Budget balancing	Truthfulness	Economic efficiency
M1	Yes	Strong	False	Yes
M2	Yes	Strong	True	Yes
M3	Yes	Weak	True	Yes

Table 1: Settlement mechanism properties

335 Three scenarios with different community supply/demand ratios are evaluated for
 336 all considered settlement mechanisms: over-supply (10:1), over-demand (1:10), and
 337 perfect balance (i.e. equal supply and demand). Each mechanism is evaluated based
 338 on the policies developed by the agents and the resulting market behaviour based
 339 on emerging equilibrium bid, ask, and settlement prices, given the same training
 340 curriculum. The goal is to find a market mechanism that follows the law of supply
 341 and demand, and is compatible with RL agent learning behavior. Such a mechanism
 342 is expected to produce the following results:

- 343 • Excess supply case: The generators compete for demand, driving ask prices low
 344 with bid prices following.
- 345 • Excess demand case: The consumers compete for supply, driving bid prices high
 346 with ask prices following.
- 347 • Equal supply and demand case: The bid and ask prices converge around the
 348 middle of the available price range.
- 349 • For all cases: The mean bid, ask, and settlement prices should have low spread.

350 A set of $n = 4$ independently learning participants is considered. Two participants
 351 with $d^i > g^i$ act as buyers, and the remaining two participants with $d^i < g^i$ act
 352 as sellers. Two of each type of participant maintain competition on both sides of
 353 the market and should prevent monopolistic behaviour. Steady-state (flat, time-
 354 invariant) energy profiles are employed for each agent, with the collective load demand
 355 and supply corresponding to the previously given ratios

$$\frac{g^{\text{LEM}}}{d^{\text{LEM}}} = \frac{\sum_i g^i}{\sum_i d^i}. \quad (26)$$

356 This setup collapses the observation space of each agent to a single point and
 357 fixes q_{bid}^i or q_{ask}^i to the residual load. This allows to further improve the purity of
 358 the experiment by learning only the price policy. From the view of a single agent,

359 this transforms the experiment into a partially observable, nonstationary multi-armed
 360 bandit, where the number of arms corresponds to the number of discrete price actions
 361 $|p|$. For each individual participant, an independent tabular Q-learning algorithm can
 362 be used, with ϵ -greedy exploration policy and learning rate α , as described in section 2.
 363 This maintains loose convergence guarantees despite the properties of the resulting
 364 environment [34]. For this experiment, α is set to 0.1, γ to 0.99, and ϵ to 0.98.
 365 Values of ϵ and α are annealed starting from episode 100 to balance exploration and
 366 convergence speed, with a multiplier of 0.98 per episode. Under this simple setup, if
 367 the agents fail to develop policies that reflect the previously mentioned criteria, the
 368 respective mechanisms will be considered infeasible for subsequent use.

369 This experiment was performed using the T-REX simulator, which is described
 370 in Appendix B. The simulations were run on a workstation with Ryzen 9 3900X
 371 processor and 32GB of 3200MHz DDR4 memory. In this specific setup, each episode
 372 took approximately 5 minutes. Detailed experimental configurations can be found on
 373 the GitHub repository of the project [35].

374 4.1.2. Results and Discussion

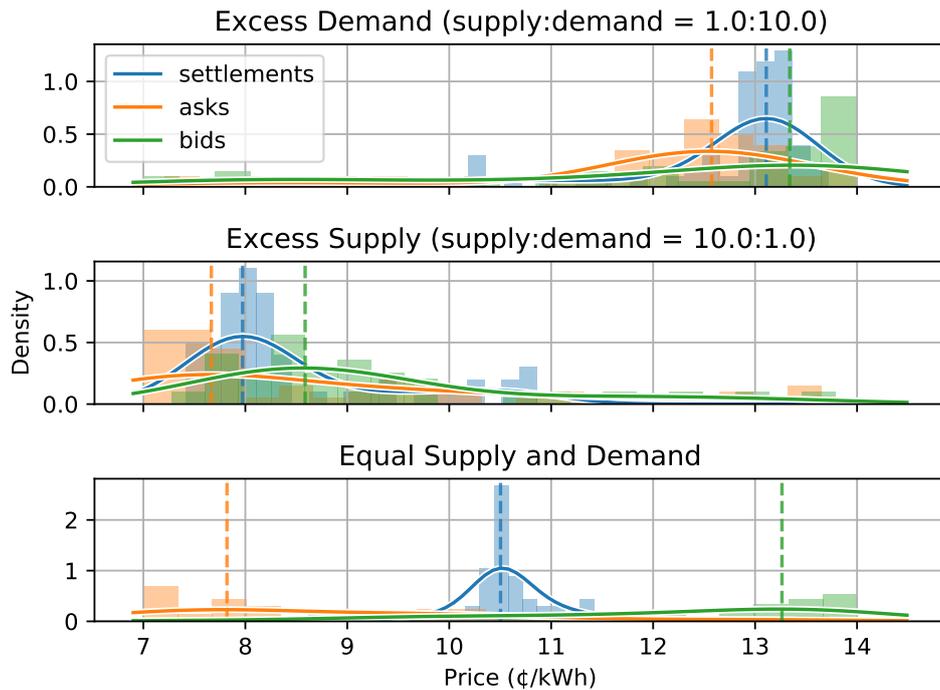


Figure 3: Validation policies for bid, ask, and resulting settlement prices for agents operating under M1 for episodes 70 to 100. Discrete policies are shown by the histograms. Probability density functions of the histograms are overlaid on top, which are approximated with the Gaussian KDE function in the scikit-learn Python package with default parameters. Means of density plots are shown by dashed lines.

375 Recall that, in ALEX, participants both determine the price signals and make
 376 energy management decisions through market interactions. Therefore, it is important

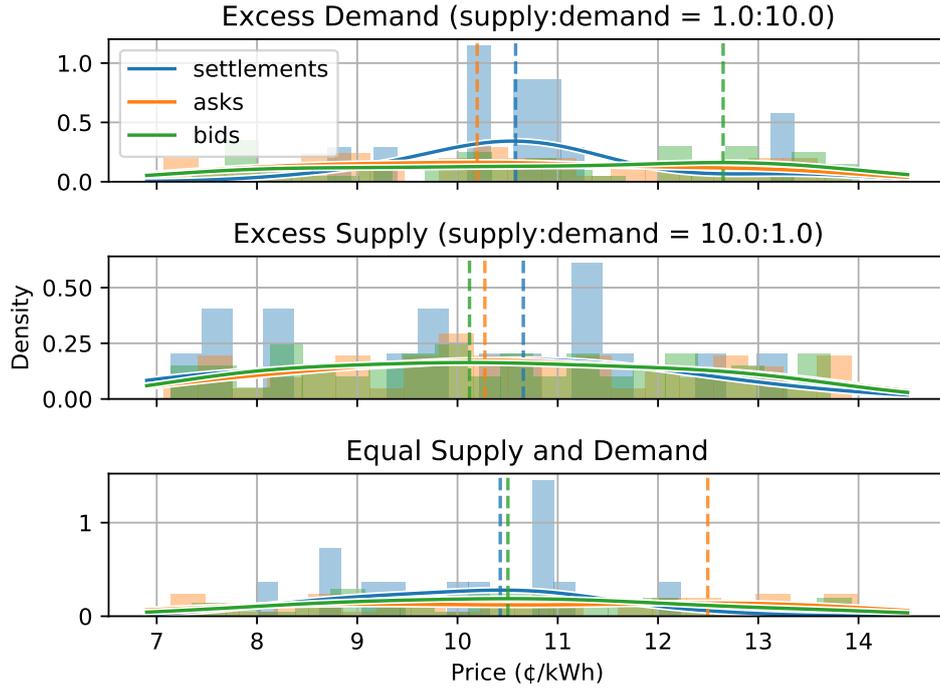


Figure 4: Validation policies for bid, ask, and resulting settlement prices for agents operating under M2 for episodes 70 to 100. Discrete policies are shown by the histograms. Probability density functions of the histograms are overlaid on top, which are approximated with the Gaussian KDE function in the scikit-learn Python package with default parameters. Means of density plots are shown by dashed lines.

377 to study the agent actions as well as the resulting settlement prices. Fig. 3 shows,
 378 for settlement mechanism M1, the policies learned for bid and ask prices, as well as
 379 the resulting settlement prices as density plots. The bid, ask, and settlement prices
 380 are, in general, closely clustered on the expected side of the price range for both
 381 unbalanced cases. However, the balanced case reveals a critical problem. While the
 382 settlement prices are concentrated in the middle of the price range (as expected),
 383 both bid and ask prices diverge near the extremities. This phenomenon results
 384 from the lack of truthfulness of M1. Agents have no incentive to submit the bid
 385 and ask prices that correspond to what they believe should be the value of energy
 386 (near the settlement price). Since M1 calculates the settlement price as the average
 387 of each pair of bid/ask, this strategy increases the chance of reaching a settlement.
 388 However, at the same time, it also increases the reward if an opponent follows a
 389 truthful strategy. This behavior is evidently the optimal strategy to employ in this
 390 scenario. However, the price divergence is problematic, especially when continuous,
 391 unbounded action spaces were used for price selection: exceedingly large bid/ask
 392 prices may cause settlement prices to become unstable. Because of this risk, M1 is
 393 disqualified.

394 Fig. 4 shows the results for settlement mechanism M2. Unlike the previous mech-
 395 anism, M2 shows no clear convergence for any case. A possible explanation is that

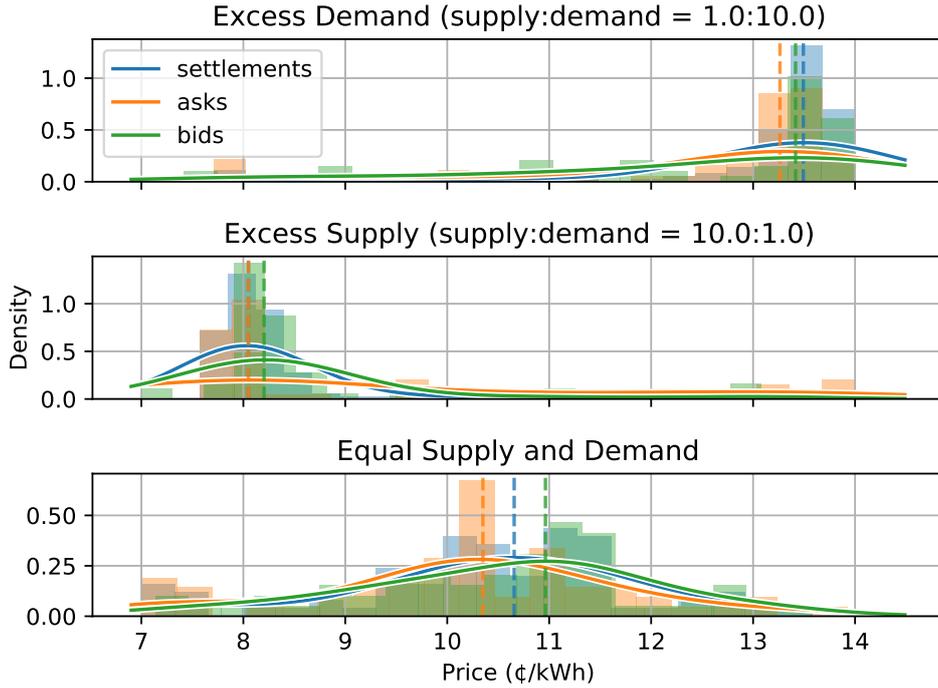


Figure 5: Validation policies for bid, ask, and resulting settlement prices for agents operating under M3 for episodes 70 to 100. Discrete policies are shown by the histograms. Probability density functions of the histograms are overlaid on top, which are approximated with the Gaussian KDE function in the scikit-learn Python package with default parameters. Means of density plots are shown by dashed lines.

396 M2 satisfies the conditions for an ideal double auction market, which is known to
 397 be impossible to practically implement according to the Myerson–Satterthwaite the-
 398 orem [33]. Another possibility is that strong budget balancing drastically decreases
 399 the number of successful settlements, which results in sparse rewards within this
 400 nonstationary environment. Therefore, despite the theoretical existence of a Nash
 401 equilibrium, agents are unlikely to discover it due to the lack of feedback. While it
 402 may be possible for M2 to converge, given a sufficiently long training time, the fact
 403 that it fails to show convergence using the same simulation parameters as M1 and
 404 M3 makes it less desirable. RL agents are often expected to update policies on data
 405 streams in real-time, which is much slower than in simulations. **The real-time equiv-**
 406 **alent of the training steps is 200 weeks, or approximately four years of one-minute**
 407 **resolution data.** However, there is still a lack of historical smart meter data. There-
 408 fore, if this system were deployed in a real environment, it would have to learn in
 409 real-time. As a result, the market mechanism M2 is disqualified as unsuitable for real
 410 world applications.

411 Fig. 5 shows the results for settlement mechanism M3. Similar to M1, the bid,
 412 ask, and settlement prices are closely clustered on the expected side of the price
 413 range for both unbalanced cases, even more closely together. However, unlike M1,
 414 the balanced case shows similar behaviour, with prices concentrated in the middle of

415 the price range. Therefore, it can be concluded that M3 has the truthfulness property.
416 The agents have little incentive to set bid/ask prices that deviate too far from the
417 settlement prices, which should closely approximate the true value of energy for each
418 ratio of supply-to-demand. Consequently, M3 is qualified for further research as it
419 satisfies the previously mentioned selection criteria.

420 The supply-to-demand ratios used in the initial experiments were quite extreme.
421 Examination of the settlement prices for more balanced ratios should provide a more
422 thorough picture of market behaviour. Therefore, the ratios 1.5:1 and 1:1.5 are added
423 to the experiments for M3. The simulations are extended by 100 episodes with an-
424 nealing as described in Section 4.1.1. As shown by the results in Fig. 6, the prices
425 settle slightly lower than the balanced case for supply-to-demand of 1.5:1, and slightly
426 higher for supply-to-demand of 1:1.5. This confirms the dominance of the law of sup-
427 ply and demand, as the settlement prices follow the supply-to-demand ratio. Further
428 experiments with more ratios of supply and demand will be performed to develop an
429 empirical model of the price behaviour.

430 In summary, the experiments show that efficient RL agent training requires weak
431 budget balancing, resulting in a stronger, denser reward signal. Truthfulness is nec-
432 essary for the emerging policies to truly reflect the law of supply and demand. In a
433 LEM with these two properties, the individual rationality of agents maximizes the
434 value exchanged between participants, guaranteeing economic efficiency as a result of
435 convergence. Even though perfect budget balancing has not been achieved, this may
436 be desirable for deployment: a small profit for the auctioneer can be used to maintain
437 the infrastructure necessary for operating the market.

438 *4.2. Economic Study*

439 *4.2.1. Experimental Design*

440 This experiment focuses on the second research question: Does the resulting mar-
441 ket behavior effectively support DR for LEM with high penetration of DER? This
442 investigation is performed by comparing the proposed approach with conventional
443 pricing schemes, such as net billing⁵ and time-of-use. As indicated by the results
444 from the previous experiment, the prices in market equilibrium are dominated by
445 the law of supply and demand (see Fig. 6). By performing additional simulations
446 with alternative proportions of supply and demand, an empirical model of the price
447 behaviour can be obtained via interpolation. Such a model can be used as a simple
448 approach to set local market prices, without implementing an actual auction-based
449 market. The supply-to-demand ratio for a local market can be derived from metering
450 data.

451 The economic study is conducted using a residential community microgrid with
452 ten participants. Due to the lack of suitable smart home data from Canada, energy
453 profiles from the openly available SunDance data set [36, 37] are used. Ten energy
454 profiles have been randomly selected to assemble the virtual community. The IDs
455 of the selected customers are as follows: 10011, 1001625, 1002714, 10068, 100703,

⁵The distinction between net billing and net metering is clarified in Appendix Appendix A.

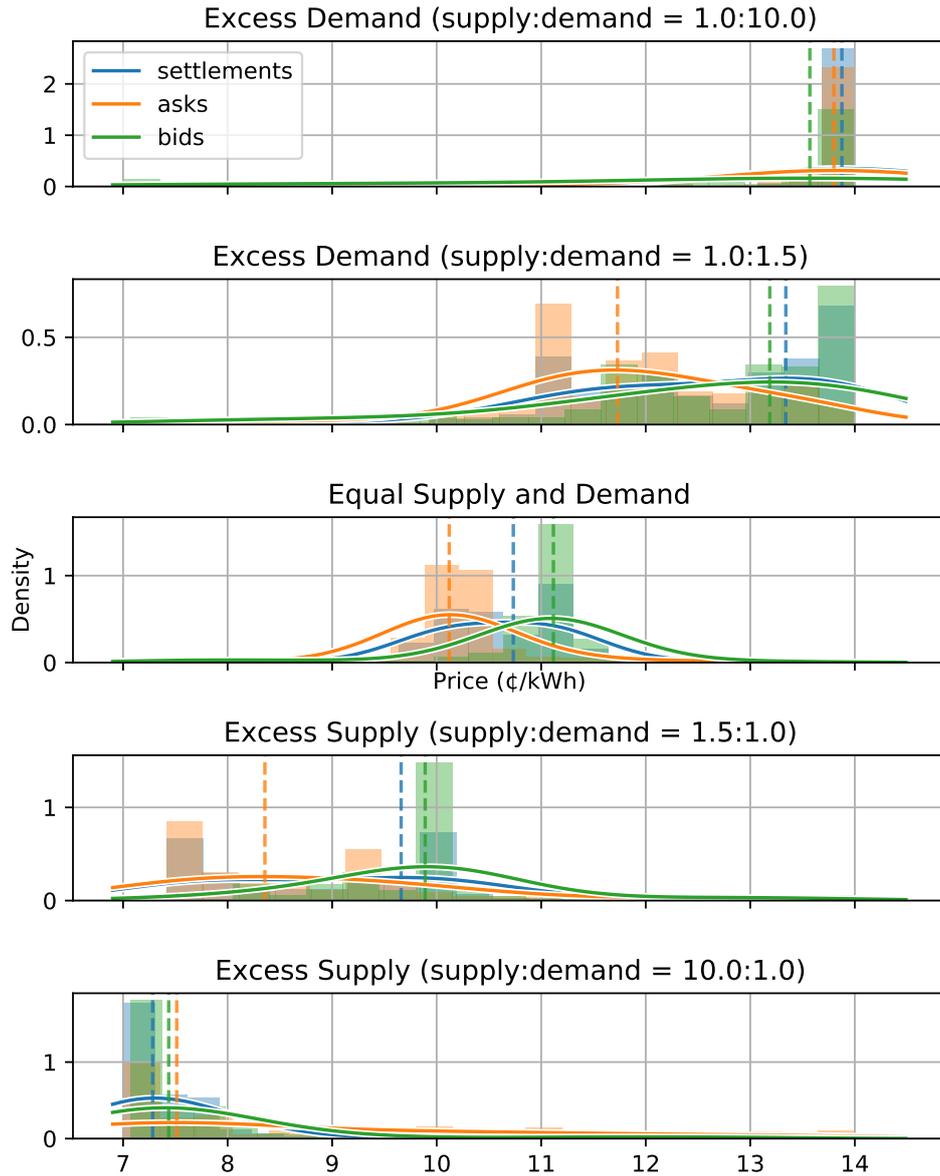


Figure 6: Validation policies for bid, ask, and resulting settlement prices for agents operating under M3 for episodes 100 to 200. Exploration factor and learning rate are annealed starting from episode 100 with a multiplier of 0.98 applied at the beginning of each episode. Discrete policies are shown by the histograms. Probability density functions of the histograms are overlaid on top, which are approximated with the Gaussian KDE function in the scikit-learn Python package with default parameters. Means of density plots are shown by dashed lines.

456 1001420, 1003173, 1001230, 100114, 100196. All participants are prosumers partici-
 457 pating in energy trading to gain economic benefits. The microgrid is assumed to be
 458 on a single bus behind a community smart meter. Similar to the previous experi-
 459 ment, no load shaping is performed. To illustrate the changes of supply-to-demand
 460 behaviour of the test community, the aggregated values of supply and demand over
 461 a single summer day (June 1, 2015) are plotted in Fig. 7.

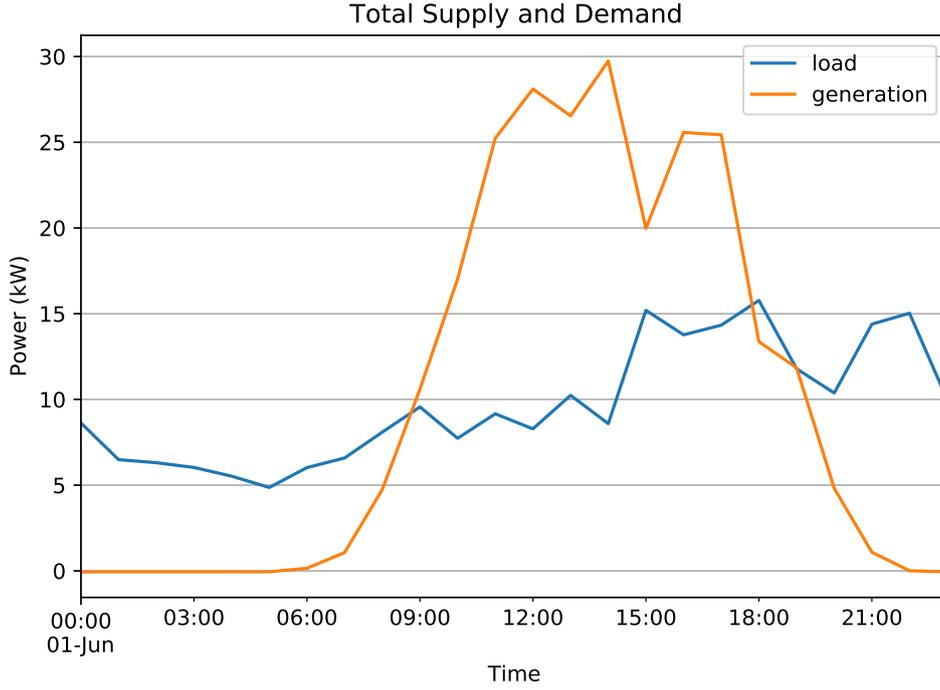


Figure 7: Total supply and demand profile of the residential community test over one summer day in June 1, 2015

462 The experiment evaluates the changes of electricity bills caused by enabling a local
 463 energy market that sets energy exchange prices based on the local supply and demand.
 464 Because the local market price is time-varying and supply-to-demand dependent, the
 465 LEM-based pricing can also be compared to a benchmark TOU pricing schedule.

466 4.2.2. Results and Discussion

467 The system-wide market model is developed using the data from the previous
 468 experiments supplemented by X additional demand ratios (values). The resulting
 469 pricing model is shown in Fig. 8. Note that in real settings, where the load demand
 470 curve and DER availability of each participant are unique, ALEX agents may develop
 471 personalized pricing schedules. The goal of this experiment is to evaluate the economic
 472 performance of an ALEX-based trading system and to compare it with common tariffs
 473 that do not use individual pricing schedules.

474 The resulting equation for the price curve is as follows:

$$P(s, d) = \begin{cases} P_{NB,load}^{grid} & P(s, d) \geq P_{NB,buy}^{grid} \\ P_{NB,gen}^{grid} & P(s, d) \leq P_{NB,sell}^{grid} \\ -0.0254 \frac{s}{d} + 0.1426 C_{H,M,L}^{TOU} & \text{if buying} \\ -0.0280 \frac{s}{d} + 0.1299 C_{H,M,L}^{TOU} & \text{if selling,} \end{cases} \quad (27)$$

475 where s is energy supply, d is energy demand, $P_{NB,load}^{grid}$ is price of electricity when

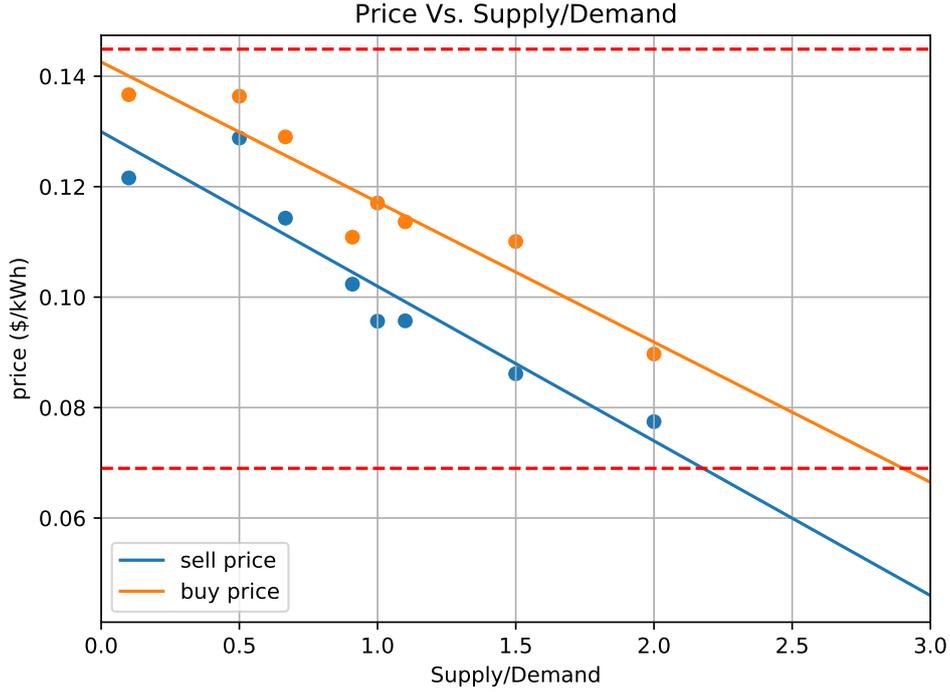


Figure 8: Pricing model developed for the test local market. The dotted lines show the price boundaries defined by (11). Linear regression between the price points leads to a well-fit, generalized mathematical model.

476 buying electricity from the grid under net billing, $P_{NB,gen}^{grid}$ is price of electricity when
 477 selling electricity to the grid under net billing, and $C_{H,M,L}^{TOU}$ are the adjustment factors
 478 when TOU is used.

479 A pricing schedule for the local market of the energy community and a selected
 480 day can be obtained by applying this model to specific energy profiles, as shown in
 481 Fig. 9 for the sample profiles from Fig. 7.

482 The internal price determined using the model corresponds well to the ratio of
 483 supply and demand of the community throughout the day. For example, at midnight,
 484 when solar generation is nil, the price for selling energy to a peer is \$0.1449, which is
 485 the same for the buyer as if purchasing energy from the grid. Later in the morning,
 486 at 7:00AM, when solar energy becomes available, the price for selling to peers lowers
 487 accordingly. At around 9:00AM, when generation is significantly higher than demand,
 488 the price for selling to peers drops to \$0.069, which is the same as selling to the grid
 489 under net billing.

490 The economic performance of all participants under this price curve is calculated
 491 and compared with net billing. The results of this comparison are shown in Fig. 10
 492 and summarized in Table 2. As a reminder, the rules of interaction between the
 493 community participants and the grid are detailed in Appendix C. In accordance with
 494 the operating principles of net billing, described in Appendix A, the entire community
 495 is placed behind a community meter, and energy exchanged directly between peers
 496 does not incur T&D fees.

Table 2: ALEX vs. Net Billing (NB)

Participant	Bill (\$)			kWh bought from Grid			kWh bought local			kWh sold to Grid			kWh sold local		
	NB	ALEX	diff.	NB	ALEX	diff.	NB	ALEX	diff.	NB	ALEX	diff.	NB	ALEX	diff.
10011	0.10	-0.28	383.2%	6.83	2.35	65.59%	N/A	4.48	100%	12.88	0	100%	N/A	N/A	12.88
1001625	0.80	0.13	83.2%	10.73	2.28	78.75%	N/A	8.45	100%	10.97	0	100%	N/A	N/A	10.97
1002714	4.34	3.26	24.8%	32.45	13.10	59.63%	N/A	19.35	100%	5.31	0	100%	N/A	N/A	5.31
10068	2.80	1.55	44.7%	22.37	4.08	81.76%	N/A	18.29	100%	6.41	0	100%	N/A	N/A	6.41
100703	2.35	1.32	44.0%	22.47	5.94	73.56%	N/A	16.53	100%	13.17	0	100%	N/A	N/A	13.17
1001420	2.21	1.47	33.6%	17.63	11.96	32.16%	N/A	5.67	100%	4.98	0	100%	N/A	N/A	4.98
1003173	7.56	5.09	32.7%	61.73	22.00	64.36%	N/A	39.73	100%	20.03	0	100%	N/A	N/A	20.03
1001230	3.55	2.69	24.2%	24.89	8.03	67.74%	N/A	16.86	100%	0.80	0	100%	N/A	N/A	0.80
100114	5.22	2.70	48.2%	49.55	12.65	74.47%	N/A	36.9	100%	28.40	0	100%	N/A	N/A	28.40
1001965	6.50	4.78	26.5%	50.78	19.07	62.45%	N/A	31.71	100%	12.45	0	100%	N/A	N/A	12.45
Total	35.43	22.70	35.9%	299.44	101.46	66.12%		198		115.40	0	100%			115.40

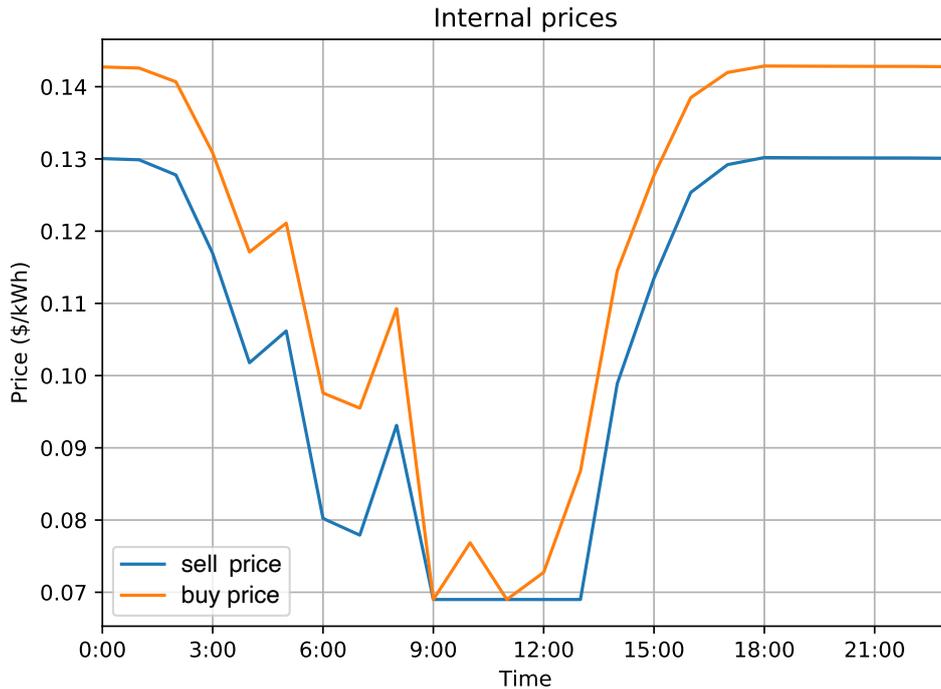


Figure 9: Internal prices of ALEX used to conduct transactions

497 The results in Table 2 show that the implementation of a community market can
 498 financially benefit the community as a whole, reducing the total community bill by
 499 35.9%. The mean and median of individual bill reductions are 74.51% and 38.8%,
 500 respectively. By putting the whole community behind-the-meter, financial benefits
 501 may even be gained by those who cannot afford the expense of acquiring and installing
 502 their own DER. This is because they have direct access to excess generation of their
 503 neighbours that may be lower priced in comparison to buying from the grid. Similarly,
 504 there is an inherent financial incentive to sell excess generation to peers first, as the
 505 profits can be higher than selling directly to the grid. In other words, this setup may
 506 further socialize the benefits of DERs.

507 As mentioned before, the local market price is time-varying and supply-to-demand
 508 dependent. This is similar to the philosophy behind the development of TOU, which
 509 uses the time-varying supply/demand of the entire grid instead of focusing on any
 510 specific area. Therefore, these two approaches are compared to quantify their relative
 511 performance. Ontario TOU is used as a benchmark in Canada and is often referenced
 512 by utility companies in jurisdictions without TOU, such as Alberta.

513 Fig. 11 displays the two pricing schedules, showing the stark contrast between their
 514 shapes. Whereas the local energy price decreases toward noon due to the increase in
 515 generation, TOU increases, which suggests that there is more load than generation
 516 during this period. While it is possible that this is true due to commercial and
 517 industrial loads, which do not exist in the local market, the fact remains that the
 518 TOU is not correlated with the actual balance of load and demand in the testing
 519 locale. While this disconnect may be a cause for the lack of participation in TOU

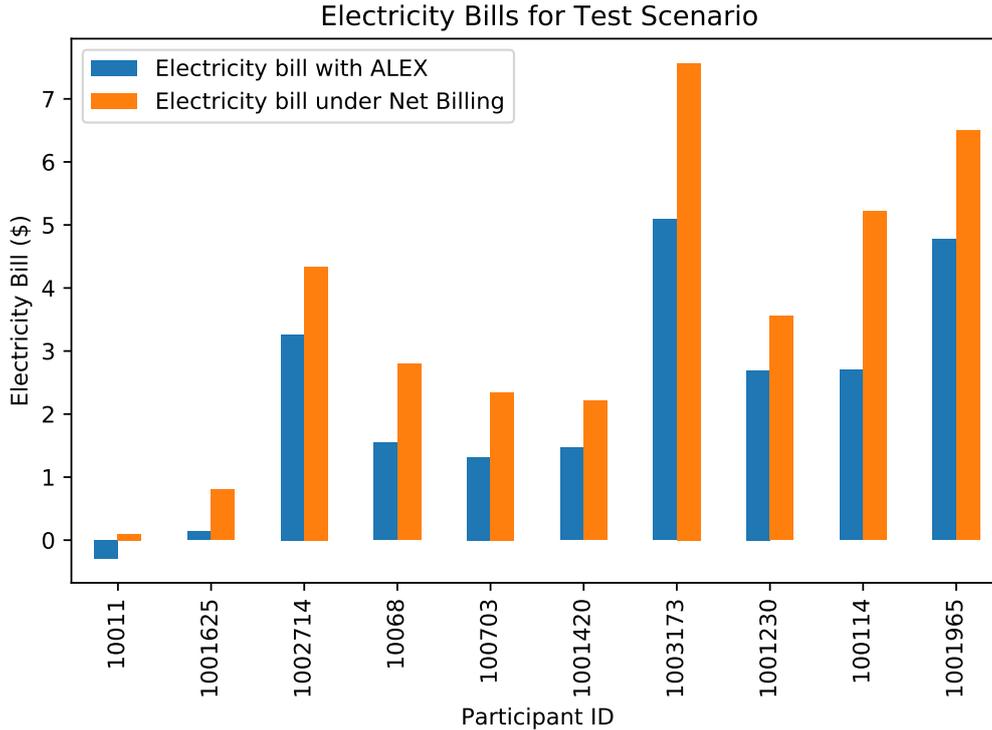


Figure 10: Electricity bill comparison between net billing and ALEX

520 mentioned in Section 2, it also suggests the need for very localized, highly relevant
 521 pricing signals to increase the efficiency of managing DERs and the overall system.
 522 ALEX is a highly scalable, distributed approach that generates highly relevant pricing
 523 signals at a low cost.

524 5. Conclusions and Future Work

525 DR techniques provide an effective means to manage DERs. As more such re-
 526 sources are installed and the energy mix becomes more complex, their management
 527 and coordination should be automated. This article explores the requirements for
 528 automating LEMs using multi-agent reinforcement learning. This exploration is fa-
 529 cilitated by ALEX, a LEM framework that can use an arbitrary closed-book, double
 530 auction settlement system. It is used to identify the market properties that drive the
 531 policies of independent Q-learning agents to follow the law of supply and demand.
 532 After establishing an appropriate market settlement mechanism, the emergent market
 533 behaviour is compared to conventional DER integration techniques.

534 The first experiment trains a group of agents with three market configurations,
 535 distinguished by their general properties. The results show that truthfulness is nec-
 536 essary for the collective policy to reflect the law of supply and demand. The second
 537 requirement, weak budget balancing, facilitates the generation of reinforcement sig-
 538 nals sufficient for the trading agents to learn.

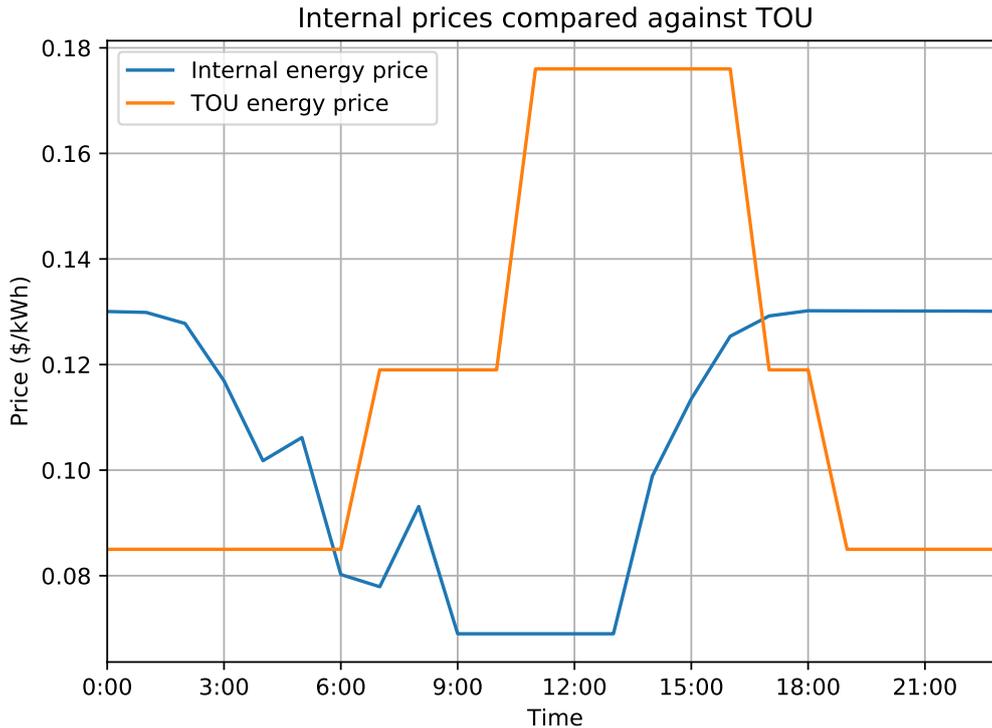


Figure 11: Local market pricing schedule compared against Ontario summer TOU prices.

539 The second experiment compares the resulting LEM behaviour with that of mar-
 540 kets based on net billing and time-of-use (TOU). Since consensus pricing in ALEX
 541 strongly reflects the law of supply and demand, the resulting price signal is signifi-
 542 cantly more responsive and relevant than TOU. In terms of economic performance
 543 in the test community, the proposed approach provides a bill reduction of 38.8%
 544 compared to net billing.

545 The findings presented in this article lay the crucial ground for future work. We
 546 plan to investigate the integration of battery energy storage systems, along with more
 547 complex RL algorithms to increase profile shaping capabilities. This may further
 548 improve economic performance, as well as grid stability. Performance comparisons
 549 with other deployed LEM approaches will be conducted to identify the most suitable
 550 automation approach for ALEX.

551 References

- 552 [1] A. Fattahi, M. Dehimi, A review of demand-side management: Reconsidering
 553 theoretical framework, *Renewable and Sustainable Energy Reviews* 80 (2017)
 554 367–379. doi:10.1016/j.rser.2017.05.207.
- 555 [2] B.-G. Kim, Y. Zhang, M. van der Schaar, J.-W. Lee, Dynamic pricing and
 556 energy consumption scheduling with reinforcement learning, *IEEE Transactions*
 557 *on Smart Grid* 7 (2016) 2187–2198. doi:10.1109/TSG.2015.2495145.

- 558 [3] R. Lu, S. Hong, X. Zhang, A dynamic pricing demand response algorithm for
559 smart grid: Reinforcement learning approach, *Applied Energy* 220 (2018) 220–
560 230. doi:10.1016/j.apenergy.2018.03.072.
- 561 [4] A. Meyabadi, M. Deihimi, A review of demand-side management:
562 Reconsidering theoretical framework, *Renewable and Sustainable En-
563 ergy Reviews* 80 (2017) 367–379. URL: [https://www.sciencedirect.com/
564 science/article/pii/S1364032117308481](https://www.sciencedirect.com/science/article/pii/S1364032117308481). doi:[https://doi.org/10.1016/
565 j.rser.2017.05.207](https://doi.org/10.1016/j.rser.2017.05.207).
- 566 [5] S. Chen, C.-C. Liu, From demand response to transactive energy: state of the
567 art, *Journal of Modern Power Systems and Clean Energy* 5 (2017) 10–19.
- 568 [6] O. Abrishambaf, F. Lezama, P. Faria, Z. Vale, Towards transactive energy sys-
569 tems: An analysis on current trends, *Energy Strategy Reviews* 26 (2019) 100418.
- 570 [7] M. Yu, R. Lu, S. Hong, A real-time decision model for industrial load manage-
571 ment in a smart grid, *Applied Energy* 183 (2016). doi:10.1016/j.apenergy.
572 2016.09.021.
- 573 [8] X. Huang, S. H. Hong, Y. Li, Hour-ahead price based energy management scheme
574 for industrial facilities, *IEEE Transactions on Industrial Informatics* 13 (2017)
575 2886–2898. doi:10.1109/TII.2017.2711648.
- 576 [9] R. de Sá Ferreira, L. A. Barroso, P. R. Lino, M. M. Carvalho, P. Valenzuela,
577 Time-of-use tariff design under uncertainty in price-elasticities of electricity de-
578 mand: A stochastic optimization approach, *IEEE Transactions on Smart Grid*
579 4 (2013) 2285–2295. doi:10.1109/TSG.2013.2241087.
- 580 [10] D. Forfia, M. Knight, R. Melton, The view from the top of the mountain: Build-
581 ing a community of practice with the gridwise transactive energy framework,
582 *IEEE Power and Energy Magazine* 14 (2016) 25–33.
- 583 [11] E. Mengelkamp, J. Diesing, C. Weinhardt, Tracing local energy markets: A
584 literature review:, *it - Information Technology* 61 (2019) 101–110. URL: [https:
585 //doi.org/10.1515/itit-2019-0016](https://doi.org/10.1515/itit-2019-0016). doi:doi:10.1515/itit-2019-0016.
- 586 [12] M. Pilz, L. Al-Fagih, Recent advances in local energy trading in the smart
587 grid based on game-theoretic approaches, *IEEE Transactions on Smart Grid* 10
588 (2019) 1363–1371. doi:10.1109/TSG.2017.2764275.
- 589 [13] M. Khorasany, Y. Mishra, G. Ledwich, Market framework for local en-
590 ergy trading: a review of potential designs and market clearing approaches,
591 *IET Generation, Transmission & Distribution* 12 (2018) 5899–5908. URL:
592 [https://ietresearch.onlinelibrary.wiley.com/doi/abs/10.1049/
593 iet-gtd.2018.5309](https://ietresearch.onlinelibrary.wiley.com/doi/abs/10.1049/iet-gtd.2018.5309). doi:<https://doi.org/10.1049/iet-gtd.2018.5309>.
594 arXiv:<https://ietresearch.onlinelibrary.wiley.com/doi/pdf/10.1049/iet-gtd.2018.5309>.

- 595 [14] B. Baker, I. Kanitscheider, T. Markov, Y. Wu, G. Powell, B. Mc-
596 Grew, I. Mordatch, Emergent tool use from multi-agent autotutorials, 2020.
597 [arXiv:1909.07528](https://arxiv.org/abs/1909.07528).
- 598 [15] E. Mengelkamp, P. Staudt, J. Ganttner, C. Weinhardt, Trading on local energy
599 markets: A comparison of market designs and bidding strategies, in: 2017 14th
600 International Conference on the European Energy Market (EEM), 2017, pp. 1–6.
- 601 [16] D. Silver, A. Huang, C. Maddison, A. Guez, L. Sifre, G. Driessche, J. Schrit-
602 twieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe,
603 J. Nham, N. Kalchbrenner, I. Sutskever, T. Lillicrap, M. Leach, K. Kavukcuoglu,
604 T. Graepel, D. Hassabis, Mastering the game of go with deep neural networks
605 and tree search, *Nature* 529 (2016) 484–489. doi:10.1038/nature16961.
- 606 [17] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung,
607 D. H. Choi, R. Powell, T. Ewalds, P. Georgiev, J. Oh, D. Horgan, M. Kroiss,
608 I. Danihelka, A. Huang, L. Sifre, T. Cai, J. P. Agapiou, M. Jaderberg, A. S.
609 Vezhnevets, R. Leblond, T. Pohlen, V. Dalibard, D. Budden, Y. Sulsky, J. Mol-
610 loy, T. L. Paine, C. Gulcehre, Z. Wang, T. Pfaff, Y. Wu, R. Ring, D. Yogatama,
611 D. Wunsch, K. McKinney, O. Smith, T. Schaul, T. Lillicrap, K. Kavukcuoglu,
612 D. Hassabis, C. Apps, D. Silver, Grandmaster level in StarCraft II using multi-
613 agent reinforcement learning, *Nature* 575 (2019). Number: 7782 Publisher: Na-
614 ture Publishing Group.
- 615 [18] OpenAI, :, C. Berner, G. Brockman, B. Chan, V. Cheung, P. Debiak, C. Denni-
616 son, D. Farhi, Q. Fischer, S. Hashme, C. Hesse, R. Józefowicz, S. Gray, C. Olsson,
617 J. Pachocki, M. Petrov, H. P. d. O. Pinto, J. Raiman, T. Salimans, J. Schlatter,
618 J. Schneider, S. Sidor, I. Sutskever, J. Tang, F. Wolski, S. Zhang, Dota 2 with
619 large scale deep reinforcement learning, 2019. [arXiv:1912.06680](https://arxiv.org/abs/1912.06680).
- 620 [19] J. Vázquez-Canteli, Z. Nagy, Reinforcement learning for demand response: A
621 review of algorithms and modeling techniques, *Applied Energy* 235 (2019) 1072–
622 89. doi:10.1016/j.apenergy.2018.11.002.
- 623 [20] H. Zang, J. Kim, Reinforcement learning based peer-to-peer energy trade
624 management using community energy storage in local energy market, *Ener-
625 gies* 14 (2021). URL: <https://www.mdpi.com/1996-1073/14/14/4131>. doi:10.
626 3390/en14144131.
- 627 [21] L. Xiao, X. Xiao, C. Dai, M. Pengy, L. Wang, H. V. Poor, Reinforcement
628 learning-based energy trading for microgrids, 2018. [arXiv:1801.06285](https://arxiv.org/abs/1801.06285).
- 629 [22] E. Foruzan, L.-K. Soh, S. Asgarpoor, Reinforcement learning approach for opti-
630 mal distributed energy management in a microgrid, *IEEE Transactions on Power
631 Systems* 33 (2018) 5749–5758. doi:10.1109/TPWRS.2018.2823641.
- 632 [23] S. Zhou, Z. Hu, W. Gu, M. Jiang, X.-P. Zhang, Artificial intelligence based smart
633 energy community management: A reinforcement learning approach, *CSEE*

- 634 Journal of Power and Energy Systems 5 (2019) 1–10. doi:10.17775/CSEEJPES.
635 2018.00840.
- 636 [24] T. Chen, W. Su, Local energy trading behavior modeling with deep reinforce-
637 ment learning, *IEEE Access* 6 (2018) 62806–62814. doi:10.1109/ACCESS.2018.
638 2876652.
- 639 [25] D. K. Gode, S. Sunder, Allocative efficiency of markets with zero-intelligence
640 traders: Market as a partial substitute for individual rationality, *Journal of*
641 *Political Economy* 101 (1993) 119–137.
- 642 [26] T. Chen, W. Su, Indirect customer-to-customer energy trading with rein-
643 forcement learning, *IEEE Transactions on Smart Grid* 10 (2019) 4338–4348.
644 doi:10.1109/TSG.2018.2857449.
- 645 [27] J.-G. Kim, B. Lee, Automatic p2p energy trading model based on reinforcement
646 learning using long short-term delayed reward, *Energies* 13 (2020). URL: <https://www.mdpi.com/1996-1073/13/20/5359>.
- 648 [28] S. Bose, E. Kremers, E. M. Mengelkamp, J. Eberbach, C. Weinhardt,
649 Reinforcement learning in local energy markets, *Energy Informatics* 4
650 (2021) 7. URL: <https://doi.org/10.1186/s42162-021-00141-z>. doi:10.
651 1186/s42162-021-00141-z.
- 652 [29] E. Mengelkamp, J. Gärttner, C. Weinhardt, Intelligent agent strategies for res-
653 idential customers in local electricity markets, in: *Proceedings of the Ninth*
654 *International Conference on Future Energy Systems, e-Energy '18*, Association
655 for Computing Machinery, New York, NY, USA, 2018, p. 97–107.
- 656 [30] I. Erev, A. E. Roth, Predicting how people play games: Reinforcement learning
657 in experimental games with unique, mixed strategy equilibria, *The American*
658 *Economic Review* 88 (1998) 848–881. URL: [http://www.jstor.org/stable/](http://www.jstor.org/stable/117009)
659 117009.
- 660 [31] J. Nicolaisen, V. Petrov, L. Tesfatsion, Market power and efficiency in a com-
661 putational electricity market with discriminatory double-auction pricing, *IEEE*
662 *Transactions on Evolutionary Computation* 5 (2001) 504–523.
- 663 [32] F. A. Oliehoek, C. Amato, *A concise introduction to decentralized POMDPs*,
664 Springer, 2016.
- 665 [33] R. B. Myerson, M. A. Satterthwaite, Efficient mechanisms for bilateral trading,
666 *Journal of Economic Theory* 29 (1983) 265–281.
- 667 [34] J. Hu, M. P. Wellman, Nash q-learning for general-sum stochastic games, *Journal*
668 *of machine learning research* 4 (2003) 1039–1069.
- 669 [35] S. Zhang, *Trex-publication-resources*, [https://github.com/sd-zhang/](https://github.com/sd-zhang/publications-resources)
670 [publications-resources](https://github.com/sd-zhang/publications-resources), 2021.

- 671 [36] D. Chen, D. Irwin, Sundance: Black-box behind-the-meter solar disaggregation,
672 in: e-Energy '17: Proceedings of the Eighth International Conference on Future
673 Energy Systems, 2017, pp. 45–55.
- 674 [37] S. Barker, A. Mishra, D. Irwin, E. Cecchet, P. Shenoy, J. Albrecht, Smart*:
675 An open data set and tools for enabling research in sustainable homes, Proc.
676 SustKDD. (2012).
- 677 [38] D. Arrachequesne, Socket.io, 2021. URL: [https://github.com/socketio/
678 socket.io](https://github.com/socketio/socket.io).
- 679 [39] EPRI, Epr distribution system simulator, 2021. URL: [https://sourceforge.
680 net/projects/electricdss/](https://sourceforge.net/projects/electricdss/).
- 681 [40] D. Krishnamurthy, Opendssdirect.py, [https://github.com/dss-extensions/
682 OpenDSSDirect.py](https://github.com/dss-extensions/OpenDSSDirect.py), 2017.
- 683 [41] L. Tesfatsion, Agent-based computational economics: growing economies from
684 the bottom up., *Artificial Life* 8 (2002) 55–82.
- 685 [42] D. Friedman, J. Rust, *The Double Auction Market: Institutions, Theories, and
686 Evidence*, Reading, Mass., 1993.
- 687 [43] D. Friedman, A simple testable model of double auction markets, *Journal of
688 Economic Behavior & Organization* (1991).

689 **Appendix A. Net Billing**

690 This appendix clarifies the distinction between net billing and net metering. In
691 certain jurisdictions, such as Alberta, Canada, the electricity market is "unbundled".
692 In simple terms, electric utilities are only in charge of building and operating the
693 infrastructure (wires), and a multitude of retailers (which cannot be the same entity
694 as the electric utility) are allowed to sell electricity to end users, with almost complete
695 freedom to set the rate of electricity. Customer bills are therefore also separated into
696 two main components: infrastructure (transmission and distribution, or T&D fees,
697 which can have a fixed component and a variable component), and energy. Under
698 net metering, any electricity that flows into the meter (loads) incurs both energy
699 and variable T&D costs, and any electricity that flows out of the meter (generation)
700 has both energy and T&D costs deducted, either as credits or cashback. Net billing
701 is the same for loads, but only the energy component is deducted for generation.
702 One way to avoid this infrastructure cost is to install both the solar panel and a
703 battery behind the meter to minimize the amount of energy flowing out of the meter.
704 The advantage of net billing is the socialized cost of infrastructure, which is more
705 evenly divided amongst all customers. In contrast, net metering tends to shift these
706 costs onto the segment of the population who cannot afford their own solar (this is
707 a commonly known and often criticized problem). The disadvantage of net billing
708 is that the return on investment (ROI) can be significantly longer due to less bill

709 deductions. From this perspective, net billing is a more fair baseline. It also provides
 710 more opportunities for community-based energy management, such as through local
 711 energy markets like ALEX.

712 Appendix B. T-REX

713 Appendix B.1. System Architecture

714 The major limitations in deploying any TE technology at a scale are communi-
 715 cation infrastructure and computational power. This is especially true for the distri-
 716 bution system, where the amount of data that needs to be collected and processed
 717 for TE is orders of magnitude greater compared to the transmission system. Further-
 718 more, the necessary infrastructure, such as SCADA, private fiber networks, voltage
 719 sensors, current sensors, etc., is typically unavailable to the distribution system and
 720 would be prohibitively expensive to retrofit.

721 The T-REX architecture is therefore designed around the least expensive way to
 722 implement and scale TE technology. This means that inexpensive, low bandwidth,
 723 long-range wireless mesh networks, such as LoRaWAN, can be used for reliable com-
 724 munications. Computing devices should also be distributed so that the total compu-
 725 tational power of the network scales with the number of TE clients. Fig. B.12 shows
 726 the simplified architecture diagram of this approach.

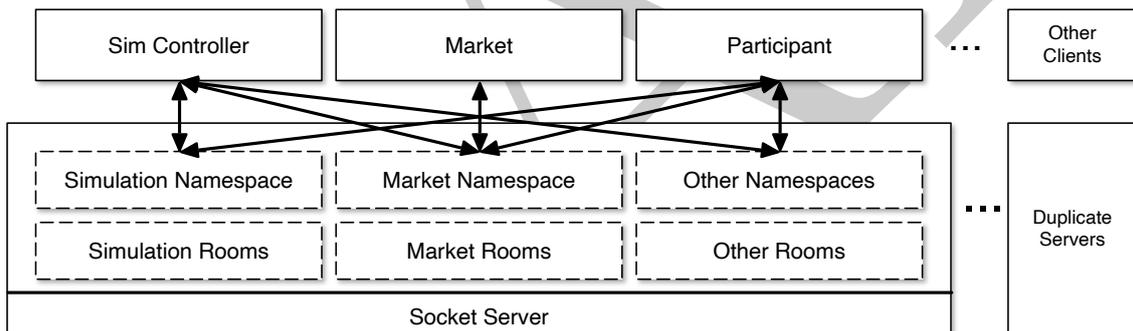


Figure B.12: Simplified T-REX V3 Architecture Diagram

727 Appendix B.2. Data Fabric

728 T-REX is built upon `socket.io` [38] as the system foundation. This guaran-
 729 tees compatibility, scalability, reliability, and deployability. When using T-REX in
 730 simulation mode, the asynchronous, highly parallel design provides true system-wide
 731 randomness and eliminates the need for pseudo-random sequence queues that are typi-
 732 cally required for contemporary TE simulators. Standard networking performance
 733 and penetration testing techniques can be readily used to evaluate the performance
 734 and cybersecurity aspects of the TE systems designed with T-REX.

735 *Appendix B.3. Clients*

736 The functional modules of T-REX are built as `socket.io` clients. As in the
737 deployment case, interaction between modules is facilitated by passing messages using
738 the `socket.io` API. Although designers are free to use payloads of any permissible
739 size, format, and endpoints, care should be taken to preserve genericity and minimize
740 bandwidth usage. There have been three main classes of clients implemented, as
741 shown in the architecture diagram and described below:

- 742 • Participant modules, which are in charge of energy trading and managing energy
743 resources that are directly accessible. Participants are, for example, households
744 and self-driving EVs.
- 745 • Non-participant modules, e.g., the TE market. The market facilitates the dis-
746 covery and exchange of energy between participants.
- 747 • Simulation-only modules, e.g., the simulation controller, or a powerflow calcula-
748 tion module. The simulation controller augments the deployment environment
749 to form a simulation model. It can also perform advanced functions such as
750 training curricula for ML applications.

751 With a few restrictions pertaining to the simulation mode, the number of modules
752 of each type is unlimited. The functions are also not restricted to the list described
753 above. For example, a traffic module can function in parallel with multiple markets
754 to guide self-driving EV participants to find optimal paths to carry passengers in
755 conjunction with charge and discharge locations to maximize profit. Other modules
756 that do not use traffic data to make decisions simply do not know about its existence.

757 *Appendix B.4. Implementing TE in T-REX*

758 TE systems can be setup in T-REX using a simple JSON configuration file. The T-
759 REX runner assembles the modules as configured and launches them as independent
760 processes on the assigned machine or machines at run time. Examples of configuration
761 files can be found on the GitHub repository [35].

762 To launch a classical TC simulation, the following modules are required:

- 763 1. A non-participant powerflow module. The built-in implementation uses OpenDSS [39]
764 and its Python API [40].
- 765 2. TC Market with a sub-module that generates prices based on the received pow-
766 erflow data.
- 767 3. Participants containing load profiles, controllable devices, and price-reactive
768 logic.

769 Fig. B.13 shows a simplified version of the sequence flow diagram of the TC co-
770 simulation implemented in T-REX. Due to the asynchronous nature of T-REX, many
771 independent asynchronous functions and parallel loops have been omitted from the
772 diagram, and only an approximation of the main flow path is shown.

773 In the same way, T-REX can also be configured to run agent-based economics
774 (ACE) [41] simulations with minimal modifications from the TC configuration. In the

775 example configurations, the only modifications are the removal of the parallel running
776 powerflow module and swapping in the appropriate market module and agent logic
777 submodules.

778 **Appendix C. Double Auction Market Design for AI**

779 *Appendix C.1. Trading Mechanism*

780 Price theory states that the price for any specific good or service is based on the
781 balance of supply and demand. In a market-based TE approach, the role of the market
782 is to efficiently facilitate the exchange of energy so that the price can appropriately
783 and accurately reflect the balance of supply and demand at the time of exchange.
784 ALEX adapts and adjusts an existing market design to fit three key considerations:

- 785 1. **Suitability for electricity grids with high penetration of DER and**
786 **RES.** This means that, from a high level perspective, a market (or a collection
787 of markets) must be able to effectively target localization and the intermittent
788 nature of RES.
- 789 2. **Technical constraints and requirements of deployment:** Data acquisition,
790 transportation, and cost must be minimized.
- 791 3. **Machine learning considerations for agents:** Related to the point above,
792 ML will play an important role in trading and managing of energy resources in
793 place of humans. For this reason, the market should be conducive to learning.
794 One way to achieve this is to compose the market with a small set of explicit
795 rules. The rules should provide a strong feedback signal, and they should be
796 flexible enough to offer large action spaces.

797 With these considerations in mind, the final market is a modified form of double
798 auctions [42][43]. The rules, explicitly implemented in the code, are described below:

- 799 1. It is assumed, for the time being, that the grid is an infinity bus and it can be
800 interacted with through net billing. We therefore adapt retail electricity prices
801 in Alberta, where buying energy from the grid costs \$0.1449/kWh, and selling
802 earns \$0.069/kWh.
- 803 2. The local market has two energy pools: one for dispatchable sources, such as
804 battery energy storage systems, and one for non-dispatchable sources, such as
805 photovoltaics. This is intended to distinguish the source of energy, and to allow
806 for the value of dispatchability to emerge.
- 807 3. Auctions settle for energy to be delivered during the one-round period from the
808 end of the current round. However, the delivery period can be parametrically
809 adjusted during run-time for future design explorations.
- 810 4. During the current round, participants submit bids and asks for energy to be
811 delivered during or beyond the next delivery period.
- 812 5. A modified double auction system is used to settle trades: bids/asks are settled
813 pairwise, with bids sorted from the highest to lowest, and asks in reverse to
814 ensure pareto equality.

- 815 6. Bid/ask quantities can be partially settled.
- 816 7. A bid/ask quantity must be an integer multiple of 1 Wh. This is in consideration
817 of future hardware integration, to allow direct use of the watt-pulse function of
818 most smart-meters.
- 819 8. During the delivery period, if a seller is in short supply, it must financially
820 compensate for the shortage at net metering prices. If batteries are available,
821 the seller has the option to compensate by discharging its batteries, for all or
822 part of the shortage during this period.
- 823 9. During the delivery period, if a buyer settled for more energy than used, the
824 buyer must still pay the seller for the unused energy at the settlement price.

825 This market design strikes a compromise between a peer-to-peer market and a
826 centralized market. By using pairwise settlements, a peer-to-peer like individualized
827 value feedback can still be provided, while the simplicity and efficiency of a centralized
828 market can be kept, especially for deployment in a small, localized region.

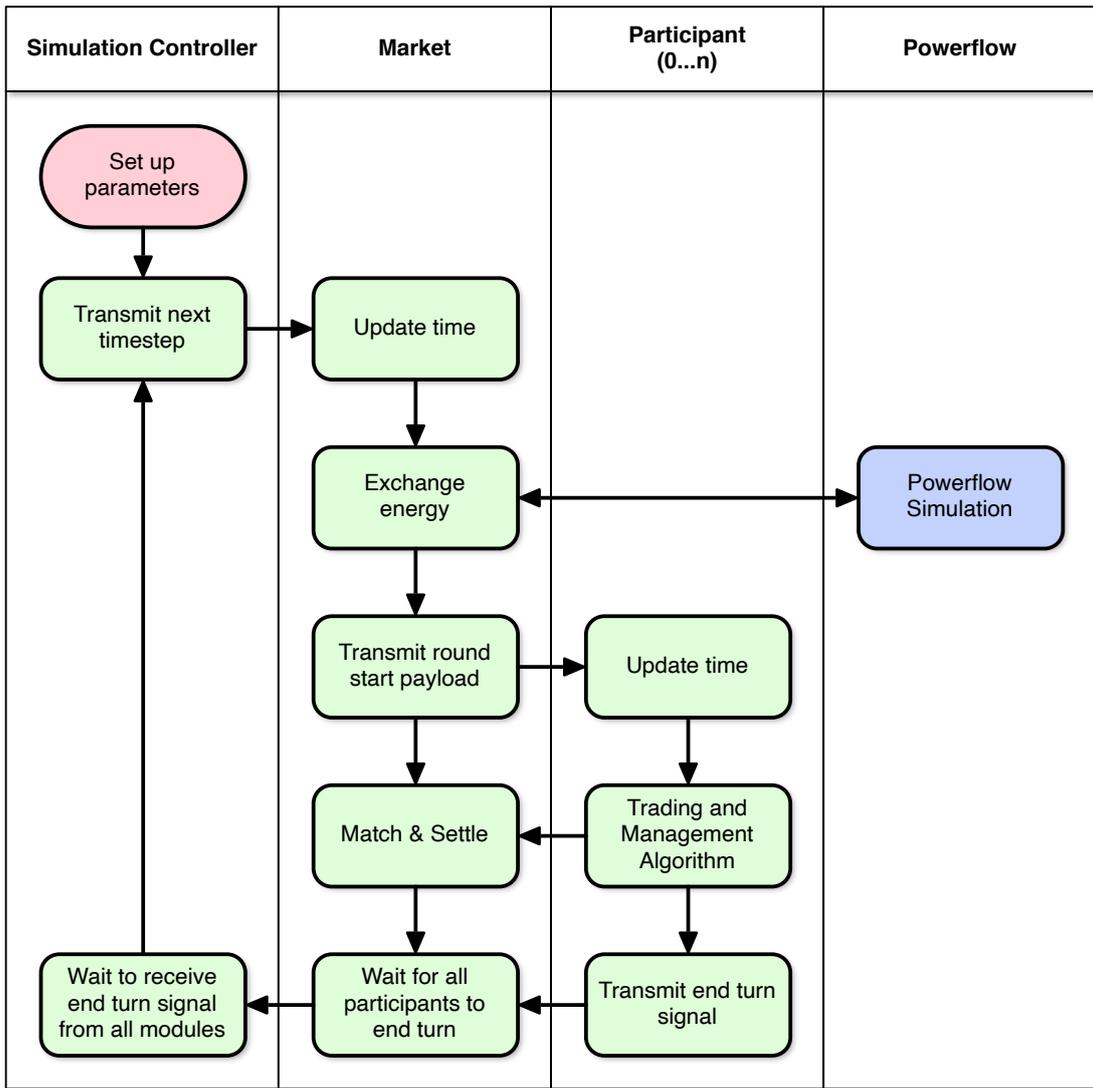
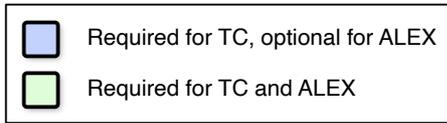


Figure B.13: Simplified swimlane diagram of TE schemes implemented in T-REX