



|                                     |   |
|-------------------------------------|---|
| <b>Title</b>                        | A Centralised Soft Actor Critic Deep Reinforcement Learning Approach to District Demand Side Management through CityLearn   |
| <b>Authors(s)</b>                   | Kathirgamanathan, Anjukan, Twardowski, Kacper, Mangina, Eleni, Finn, Donal  |
| <b>Publication date</b>             | 2020-11-17  |
| <b>Publication information</b>      | Kathirgamanathan, Anjukan, Kacper Twardowski, Eleni Mangina, and Donal Finn. "A Centralised Soft Actor Critic Deep Reinforcement Learning Approach to District Demand Side Management through CityLearn." ACM, November 17, 2020. <a href="https://doi.org/10.1145/3427773.3427869">https://doi.org/10.1145/3427773.3427869</a> . |
| <b>Conference details</b>           | The 1st International Workshop on Reinforcement Learning for Energy Management in Buildings & Cities (RLEM 2020), New York, United States of America, 17 November 2020  |
| <b>Publisher</b>                    | ACM   |
| <b>Item record/more information</b> | <a href="http://hdl.handle.net/10197/11853">http://hdl.handle.net/10197/11853</a>   |
| <b>Publisher's version (DOI)</b>    | 10.1145/3427773.3427869   |

Downloaded 2026-05-01 23:34:34

The UCD community has made this article openly available. Please share how this access benefits you. Your story matters! (@ucd\_oa)



© Some rights reserved. For more information

# A Centralised Soft Actor Critic Deep Reinforcement Learning Approach to District Demand Side Management through CityLearn

ANJUKAN KATHIRGAMANATHAN\*, School of Mechanical and Materials Engineering, University College Dublin, Ireland

KACPER TWARDOWSKI and ELENİ MANGINA, School of Computer Science, University College Dublin

DONAL P. FINN, School of Mechanical and Materials Engineering, University College Dublin, Ireland

Reinforcement learning is a promising model-free and adaptive controller for demand side management, as part of the future smart grid, at the district level. This paper presents the results of the algorithm that was submitted for the CityLearn Challenge, which was hosted in early 2020 with the aim of designing and tuning a reinforcement learning agent to flatten and smooth the aggregated curve of electrical demand of a district of diverse buildings. The proposed solution secured second place in the challenge using a centralised ‘Soft Actor Critic’ deep reinforcement learning agent that was able to handle continuous action spaces. The controller was able to achieve an averaged score of 0.967 on the challenge dataset comprising of different buildings and climates. This highlights the potential application of deep reinforcement learning as a plug-and-play style controller, that is capable of handling different climates and a heterogenous building stock, for district demand side management of buildings.

CCS Concepts: • **Applied computing** → **Engineering**; • **Computing methodologies** → **Multi-agent reinforcement learning**.

Additional Key Words and Phrases: Deep Reinforcement Learning, Smart Grid, Demand Side Management

## ACM Reference Format:

Anjukan Kathirgamanathan, Kacper Twardowski, Eleni Mangina, and Donal P. Finn. 2020. A Centralised Soft Actor Critic Deep Reinforcement Learning Approach to District Demand Side Management through CityLearn. In *The 1st International Workshop on Reinforcement Learning for Energy Management in Buildings & Cities (RLEM’20)*, November 17, 2020, Virtual Event, Japan. ACM, 7 pages. <https://doi.org/10.1145/3427773.3427869>

## 1 INTRODUCTION

Buildings are complex systems influenced by changing weather, occupancy, schedules and in a demand response (DR) context - grid signals. Capturing these dynamics in a physics-based simulation model capable of facilitating DR is highly challenging. This is due to highly non-linear behaviour of the thermal-dynamics and the fact that no one building is identical to another in the highly heterogenous building stock [3]. Reinforcement learning (RL) is a model-free algorithm that learns from historical and real-time data and has shown promise in recent research applied to building energy management problems [8, 9]. Given the novelty associated with RL in this domain, the behaviour of multiple energy consuming agents (i.e., buildings), subject to demand-dependent grid signals, is an area that is not well understood [8]. The CityLearn project (<https://github.com/intelligent-environments-lab/CityLearn>) is an OpenAI Gym environment [1], which aims to facilitate the implementation of RL agents in a multi-agent DR context for a diverse group of buildings [7]. The main objective of CityLearn is to facilitate and standardize the evaluation and comparison of different RL

---

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2020 Copyright held by the owner/author(s).

Manuscript submitted to ACM

agents and algorithms. The CityLearn Challenge was organised virtually and ran from January to July 2020. It invited participants to design, develop and tune a RL agent to flatten and smooth the aggregated curve of electrical demand for a district comprising of 9 diverse buildings. The current paper presents the results from a centralised "Soft Actor Critic" deep RL based algorithm that was submitted to the challenge.

## 2 RELATED WORK

The review of RL for DR by Vázquez-Canteli and Nagy [8] shows its promising potential as a model-free technique, mitigating the need to develop physics-based control-oriented models, and capable of dealing with the heterogenous nature of the building stock. The review found that most studies to date focus on single building systems with demand-independent electricity prices. Focusing on deep RL, which has gained significant interest and traction in recent years, e.g., using Deep Q Networks [5], such approaches have often been limited to discrete and low-dimensional action spaces [4]. There is a research gap in the application of deep reinforcement learning to problems with continuous action spaces in the building energy management domain.

The Soft Actor-Critic (SAC) algorithm, an off-policy maximum entropy actor-critic algorithm, as first proposed by Haarnoja et al [2] in 2018, is one of the algorithms that is capable of operating over continuous action spaces. At their core, actor-critic methods are a type of policy gradient method, which have separate memory structures to explicitly represent the policy [6]. The policy structure is known as the actor and the estimated value function is known as the critic. The actor selects the actions whereas the critic evaluates the actions made by the actor. The reader is referred to Sutton [6] for a more detailed explanation of Actor-Critic methods. Haarnoja et al. [2] suggest that the SAC algorithm provides for both sample-efficient learning and stability and hence extends readily to complex, high-dimensional tasks. They found the SAC algorithm showed substantial improvement in both performance and sample efficiency over both off-policy and on-policy prior methods. This current research investigates the suitability of the SAC algorithm for tackling the district DSM problem utilising CityLearn.

## 3 METHODS

### 3.1 CityLearn Challenge

The CityLearn challenge used a multi-objective cost function of five equally weighted metrics applied to an entire district of nine buildings (as outlined in Table 1). These are described below:

- (1) Peak electricity demand (for the evaluation period of 1 year)
- (2) Average daily electricity peak demand (daily peak demand of the district averaged over the evaluation period)
- (3) Ramping (a measure of how much the district electricity consumption changes from one timestep to the next)
- (4) 1 - Load factor (the average monthly electricity demand divided by its maximum peak)
- (5) Net electricity consumption of the district over the evaluation period

The multi-objective cost function is normalised by a baseline cost obtained from the performance of a predefined manually tuned Rule-Based Controller (RBC). This implied that a cost function of less than 1 resulted in a better performance than the RBC. This RBC controller charges cooling (and domestic hot water (DHW) if available) during the night and discharges during the day based only on the hour of the day. The adaptive potential of RL to deal with different environments, rather than solely the one it was trained on, was tested through evaluating the controller on different datasets. Participants used the design dataset to implement their RL agent (including design, tuning and pre-training) and could test their agent on the evaluation dataset and receive feedback through the cost function based

on the generalisation results of their agents. Each dataset contains year-long hourly information about the cooling and DHW demand of the building, electricity consumed by appliances, solar power generation, as well as weather data and other variables. The evaluation dataset featured different buildings from different cities than the design dataset, albeit within the same climate zones (see Table 2). The challenge dataset is different from both the design and evaluation datasets featuring different buildings and climates.

Table 1. Buildings and Descriptions in CityLearn Challenge District (Design Dataset)

| Building Number | Type        | Type Details         | Cooling Storage <sup>1</sup> | DHW Storage <sup>1</sup> | PV (kW) |
|-----------------|-------------|----------------------|------------------------------|--------------------------|---------|
| 1               | Commercial  | Medium Office        | 3                            | 3                        | 120     |
| 2               | Commercial  | Fast-food Restaurant | 3                            | 3                        | N/A     |
| 3               | Commercial  | Standalone Retail    | 3                            | N/A                      | N/A     |
| 4               | Commercial  | Strip Mall Retail    | 3                            | N/A                      | 40      |
| 5               | Residential | Medium Multi-family  | 3                            | 3                        | 25      |
| 6               | Residential | Medium Multi-family  | 3                            | 3                        | 20      |
| 7               | Residential | Medium Multi-family  | 3                            | 3                        | N/A     |
| 8               | Residential | Medium Multi-family  | 3                            | 3                        | N/A     |
| 9               | Residential | Medium Multi-family  | 3                            | 3                        | N/A     |

<sup>1</sup>The storage capacity is the non-dimensional scaling factor given above multiplied by the building maximum cooling or DHW demand.

Table 2. Climate Zones for the Different Datasets in the CityLearn Challenge

| CityLearn Climate Zone    | ASHRAE Identifier | Description | City        |
|---------------------------|-------------------|-------------|-------------|
| <b>Design Dataset</b>     |                   |             |             |
| 1                         | 2A                | Hot-Humid   | New Orleans |
| 2                         | 3A                | Warm-Humid  | Atlanta     |
| 3                         | 4A                | Mixed-Humid | Nashville   |
| 4                         | 5A                | Cold-Humid  | Chicago     |
| <b>Evaluation Dataset</b> |                   |             |             |
| 1                         | 2A                | Hot-Humid   | Orlando     |
| 2                         | 3A                | Warm-Humid  | Dallas      |
| 3                         | 4A                | Mixed-Humid | Kansas City |
| 4                         | 5A                | Cold-Humid  | Omaha       |

### 3.2 State Space and Hyperparameters

The state space is what the RL agent observes for each control step. The CityLearn environment allows for a total of 27 observations per building that may be passed to the agent. The final state design used in the submission was determined by utilising a combination of expert assessment and trial and error and is outlined in Table 3. A centralised solution is

presented here with one agent, which has complete oversight of all nine buildings. The SAC RL algorithm used has several key hyperparameters and the values used in the submission are detailed in Table 4. Note that some parameters were modified in the deployment (evaluation) phase to allow the agent to adapt to new environments (such as a new climate), whilst also retaining the initial weights of the pre-trained agent.

Table 3. State Space used for Centralised Agent Case

| State Variable           | Description   |
|--------------------------|---|
| month                    | month of timestep (1-12)                                    |
| day                      | day of timestep (1-7)                                       |
| hour                     | hour of the day (1-24)                                      |
| $t_{out}$                | Outside drybulb temperature (C)                             |
| direct_solar_rad         | Direct solar radiation ( $W/m^2$ )                          |
| non_shiftable_load $_n$  | Non-shiftable electricity load of Building $n$              |
| solar_gen $_n$           | Solar generation of Building $n$ (if PV present)            |
| cooling_storage_soc $_n$ | State of charge of cooling storage of Building $n$          |
| dhw_storage_soc $_n$     | State of charge of DHW storage of Building $n$ (if present) |

Table 4. Hyperparameters for Training and Evaluation

| Symbol   | Description                  | Training Value     | Evaluation Value   |
|----------|------------------------------|--------------------|--------------------|
|          | Replay buffer size           | $2 \times 10^6$    | $2 \times 10^6$    |
|          | Minibatch size               | 1024               | 64                 |
| $\gamma$ | Discount factor              | 0.9                | 0.9                |
| $\alpha$ | Reward temperate parameter   | 0.2                | 0.2                |
|          | Update interval              | 168                | 168                |
|          | Learning rate                | $5 \times 10^{-4}$ | $1 \times 10^{-4}$ |
| $\tau$   | Target smoothing coefficient | $3 \times 10^{-3}$ | $3 \times 10^{-3}$ |
|          | Hidden layer size            | 256                | 256                |

### 3.3 Reward Function

The reward function was designed based on a virtual price signal (penalising peak consumption) together with manual reward shaping to incentivise charging during the night and discharging during the day. The reward (R) function used is shown in Eq. 1 where  $\beta$  is a weighting coefficient (with a value of 0.005), N is the number of buildings in the district,  $e_{total}$  is the total district electricity consumption and  $e_i$  is the building  $i$  electricity consumption. The final reward value was also scaled and clipped to be in the range of -1 to 1.

$$R = \sum_{j=1}^N (\beta * e_{total} * e_i) + R_{night} + R_{day} \quad (1)$$

Table 5. Evaluation Results on the Design &amp; Challenge Dataset

*Note: The values presented below are relative to the baseline predefined RBC controller which has a score of 1. The lower the score, the better the performance of the RL agent.*

| Climate Zone             | Ramping | 1-Load Factor | Avg. Daily Peak | Peak Demand | Net Elec. Consumption | Avg. Score |              |
|--------------------------|---------|---------------|-----------------|-------------|-----------------------|------------|--------------|
| <b>Design Dataset</b>    |         |               |                 |             |                       |            |              |
| 1                        | 0.735   | 0.881         | 0.849           | 0.986       | 1.014                 | 0.893      |              |
| 2                        | 0.777   | 1.017         | 0.940           | 1.187       | 1.018                 | 0.988      |              |
| 3                        | 0.810   | 0.983         | 0.985           | 1.077       | 1.019                 | 0.975      |              |
| 4                        | 0.789   | 0.983         | 0.959           | 1.004       | 1.014                 | 0.950      |              |
|                          |         |               |                 |             |                       | Avg. Score | 0.952        |
| <b>Challenge Dataset</b> |         |               |                 |             |                       |            |              |
| 1                        | 0.779   | 1.014         | 0.982           | 1.131       | 1.015                 | 0.984      |              |
| 2                        | 0.780   | 0.980         | 0.959           | 0.999       | 1.013                 | 0.946      |              |
| 3                        | 0.812   | 0.960         | 0.939           | 1.083       | 1.018                 | 0.962      |              |
| 4                        | 0.860   | 0.996         | 0.991           | 1.013       | 1.017                 | 0.976      |              |
|                          |         |               |                 |             |                       | Avg. Score | <b>0.967</b> |

where:

$$R_{night} = \begin{cases} 1000, & \text{if } 10pm \leq hour \leq 12pm \text{ AND } mean(actions) > 0.1 \\ -1000, & \text{if } 10pm \leq hour \leq 12pm \text{ AND } mean(actions) < 0 \\ 0, & \text{otherwise} \end{cases}$$

$$R_{day} = \begin{cases} -1000, & \text{if } 12pm \leq hour \leq 08pm \text{ AND } mean(actions) > 0 \\ 0, & \text{otherwise} \end{cases}$$

#### 4 RESULTS

The centralised SAC RL Agent shows promising performance applied to the district DSM problem. Within 10 episodes of training on the dataset for climate zone 1 (see Fig. 1), the agent realises an improved multi-objective cost function (as defined in Section 3.1) as compared to the RBC baseline (manually predefined controller). Note that the cost function is computed relative to the RBC baseline, i.e., a cost function of less than 1 is considered to be an improved performance over the baseline and a cost function of greater than 1 is considered to be a poorer performance compared to the baseline. When this pre-trained agent is evaluated (i.e., deployed) for the same climate zone, it produces an improvement of 10.7% (a score of 0.893) over the RBC baseline (see Figure 2 for the district electricity consumption profile and Table 5 for the scores). This table also shows the ability of the RL agent to generalise and adapt to new climatic conditions, and although they suffered a performance drop, still managed to outperform the manually tuned RBC.

The Challenge dataset scores showed that the agent was further generalisable to unseen data (see Table 5) featuring both different building properties and climates. Overall, an average improvement of 3.3% was seen over the manually tuned RBC over the four different climate zones tested. Whilst these improvements are modest, they are promising given the adaptability over a range of buildings and climates and the limited information required for the state space.

#### 5 CONCLUSIONS AND FURTHER WORK

In this paper, a centralised ‘Soft Actor Critic’ reinforcement learning agent, capable of handling continuous action spaces, is proposed for the district demand side management problem and the performance of the agent applied to the

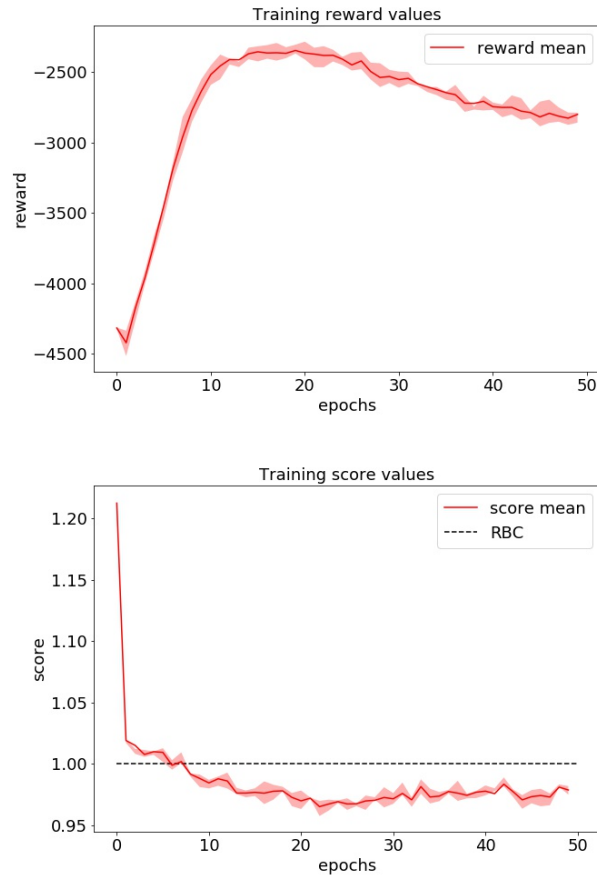


Fig. 1. Training results for Climate Zone 1 (reward and cost objectives as function of training episodes).

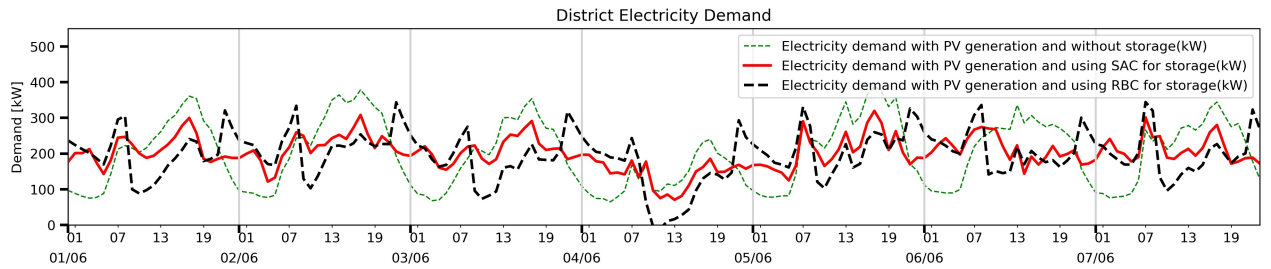


Fig. 2. A comparison of district electricity consumption for (i) no load shifting, (ii) predefined RBC and (iii) SAC RL.

CityLearn challenge is outlined. The agent was able to secure second place in the competition achieving an average score of 0.967 over the challenge dataset featuring different buildings and climates when compared to the reference manually tuned rule-based controller. This highlights the potential of deep reinforcement learning as a plug-and-play

style controller for district level demand side management of buildings. Limitations include the manual reward shaping applied which perhaps limits the generalisation ability of the RL agent to districts with significantly different demand profiles. Given the centralised agent with oversight of all buildings, it is not known how the computational requirements and performance would scale over a larger number buildings. A further limitation of CityLearn is that the cooling load is precomputed and hence currently does not support thermal comfort considerations and utilisation of the passive thermal mass for load shifting. Future work aims to further the robustness of the RL agent through reducing the amount of manual reward shaping applied and testing the performance of the algorithm for different hyperparameters. The addition of further energy systems such as batteries and electric vehicles to CityLearn will also be considered.

## ACKNOWLEDGMENTS

The authors gratefully acknowledge that their contribution emanated from research supported by Science Foundation Ireland under the SFI Strategic Partnership Programme Grant Number SFI/15/SPP/E3125.

## REFERENCES

- [1] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. 2016. OpenAI Gym. arXiv:arXiv:1606.01540
- [2] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *35th International Conference on Machine Learning, ICML 2018*, Vol. 5. 2976–2989. arXiv:arXiv:1801.01290v2
- [3] Anjukan Kathirgamanathan, Mattia De Rosa, Eleni Mangina, and Donal P. Finn. 2020. Data-driven Predictive Control for Unlocking Building Energy Flexibility: A Review. *Renewable and Sustainable Energy Reviews* 135, January 2021 (2020), 110120. <https://doi.org/10.1016/j.rser.2020.110120>
- [4] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. 2016. Continuous control with deep reinforcement learning. In *4th International Conference on Learning Representations, ICLR 2016*. arXiv:1509.02971
- [5] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. 2015. Human-level control through deep reinforcement learning. *Nature* 518, 7540 (2015), 529–533. <https://doi.org/10.1038/nature14236>
- [6] Andrew G Sutton, Richard S.; Barton. 2014. *Reinforcement Learning: An Introduction* (second ed.). MIT Press, Cambridge, Massachusetts.
- [7] José R. Vázquez-Canteli, Jérôme Kämpf, Gregor Henze, and Zoltan Nagy. 2019. CityLearn v1.0: An OpenAI gym environment for demand response with deep reinforcement learning. *BuildSys 2019 - Proceedings of the 6th ACM International Conference on Systems for Energy-Efficient Buildings, Cities, and Transportation* (2019), 356–357. <https://doi.org/10.1145/3360322.3360998>
- [8] José R. Vázquez-Canteli and Zoltán Nagy. 2019. Reinforcement learning for demand response: A review of algorithms and modeling techniques. *Applied Energy* 235, November 2018 (2019), 1072–1089. <https://doi.org/10.1016/j.apenergy.2018.11.002>
- [9] Zhe Wang and Tianzhen Hong. 2020. Reinforcement Learning for Building Controls: The problem, opportunities and challenges. *Applied Energy* 269, 1 (2020), 300. <https://doi.org/10.1016/j.apenergy.2020.115036>