

# Multi-agent reinforcement learning for renewable integration in the electric power grid

**Vincent Mai<sup>1</sup>, Tianyu Zhang<sup>1</sup> & Antoine Lesage-Landry<sup>2</sup>**

<sup>1</sup>Mila, <sup>2</sup>Polytechnique Montréal & GERAD

Tackling Climate Change with Machine Learning: workshop at NeurIPS 2021

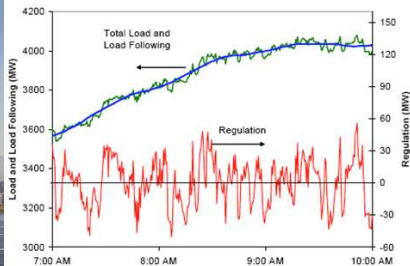
Tuesday December 14<sup>th</sup>, 2021

# Renewable integration & challenges



(a) Uncertain, **intermittent**, fast-ramping **renewables** Source:

NRDC/Vanja Terzic/iStock



(b) Sustained need for **balancing generation & demand** Source: Oak Ridge

National Lab

Figure 1: Paradigm shift in electric grid operations

- **Frequency regulation**: load balancing on short timescale;
- **Stability** of the grid;
- **Demand response of thermostatically controlled loads**
  - ▶ **modulate power consumption** of air conditioners/heating to mitigate renewable intermittency.
- Alternatives: fuel-burning power plant or expensive batteries.

# Multi-agent reinforcement learning-based demand response

**Core challenges** in frequency regulation with thermostatic loads:

① **Uncertainty**:

- limited measurements or feedback;
- exogenous factors to the power systems, e.g., weather.

② **Real-time** decision-making:

- requires power adjustment every few seconds.

③ **Dynamic** load constraints:

- temperature preferences & lock out of air conditioners.

④ **Scalability** to & **cooperativity** in multi-agent settings:

- coordination of many loads → consumption meets objective.

**Multi-agent PPO** to **coordinate the power consumption** of residential households and provide **frequency regulation**.

① policy → adaptive;

② policy → only requires evaluation to obtain decision;

③ Bellman equation in RL → dynamic environments;

④ centralized training & decentralized execution → scalability.

# Problem formulation

The demand response process is set as an POMDP. At time  $t$ :

- ① Each TCL  $i$  observes its **state**  $\mathbf{X}_{i,t}$ :
  - current, target and outdoors temperatures
  - AC power and lockdown status
  - regulation signal and aggregated consumption
  - time and date
  - thermal characteristics of the TCL
  - communication from neighbours
- ② Each TCL takes an **action**  $u_{t,i}$ : turn the TCL ON or OFF.
- ③ Each TCL receives a **reward**  $r_{t,i}$  based on:
  - the tracking of the temperature target (different for each TCL)
  - the tracking of the regulation signal by the aggregated power consumption (same for all TCLs)

# Multi-agent environment

The multi-agent environment is based on OpenAI Gym and Ray.

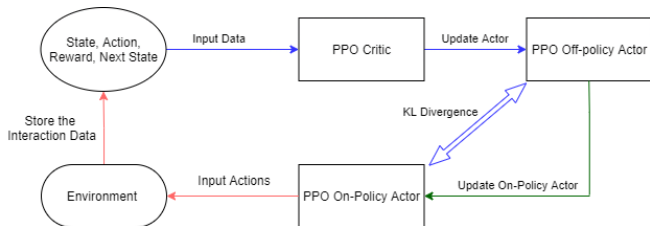
- **multiple TCLs**: from 1 to +1000
- second order **thermal model** inspired by GridLAB-D
- outdoors temperature using **real data**
- noisy **regulation signal**: first artificial, then real data.



Figure 2: Rendering of the environment.

# Multi-agent reinforcement learning

- Proximal Policy Optimization (PPO)
- Off-policy policy gradient
- Stability and performance improvements by adding constraints between the on-policy and the off-policy actors.



**Figure 3:** 2-stage PPO training scheme: 1) data gathering and 2) agent training.

# Multi-agent reinforcement learning

- **Parameter sharing** with centralized training and decentralized execution
- **Hindsight experience replay**: a multi-goal, single agent problem
- **Curriculum Learning**: design sub-goals for the agents to boost training

# Conclusion & collaboration

To tackle the **uncertainty** and **dynamic** aspects of thermostatic loads in DR, we formulate a **multi-agent reinforcement learning** approach for **frequency regulation**.

Looking for **collaborators** for:

- Building & smart grid **modelling** expertise;
- Building **simulators** (general, AC/heating or water heaters);



# Future work

- Safety: network-aware demand response;
- Simulation to real implementation (sim-to-real);
- Extension to other flexible loads (e.g. EVs.);
- Experimental case study.

**Vincent Mai**

vincent.mai@umontreal.ca

**Tianyu Zhang**

tianyu.zhang@mila.quebec

**Antoine Lesage-Landry**

antoine.lesage-landry@polymtl.ca

This work was funded by the Institute for Data Valorization (IVADO), by the National Sciences and Engineering Research Council of Canada, as well as Microsoft and Samsung.



# Appendix