



BUILDING CONTROL USING DEEP REINFORCEMENT LEARNING

Tanmay Ambadkar

Rosina Adhikari

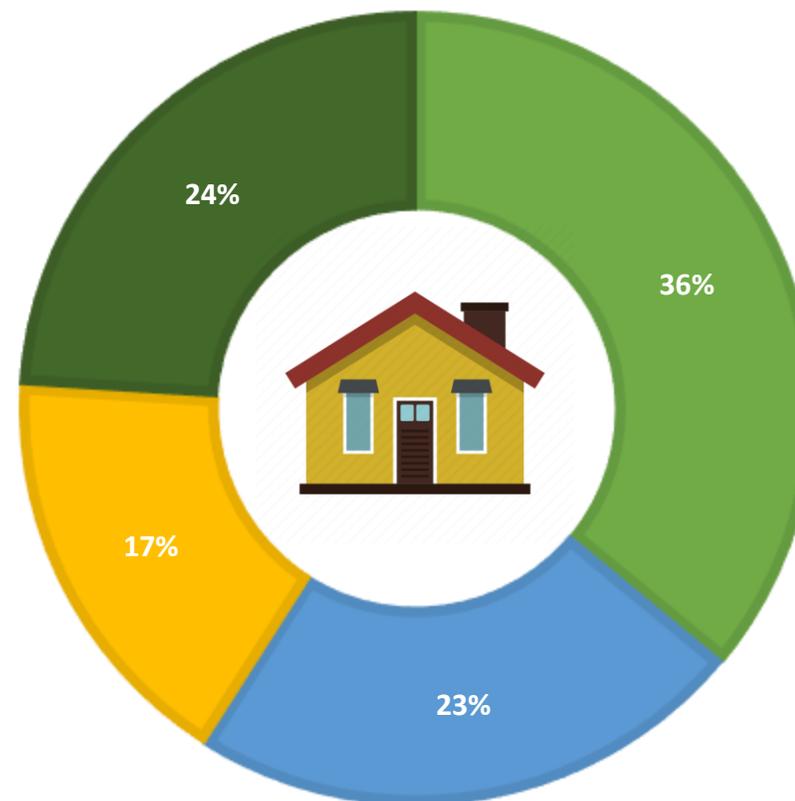
Background

According to the U.S. Energy Information Administration (EIA), in 2020, the residential and commercial **building** sectors represented about **35%** of the total energy consumed in the United States.

Space **heating** and space **cooling** account for around **60%** of total energy consumed in a building.

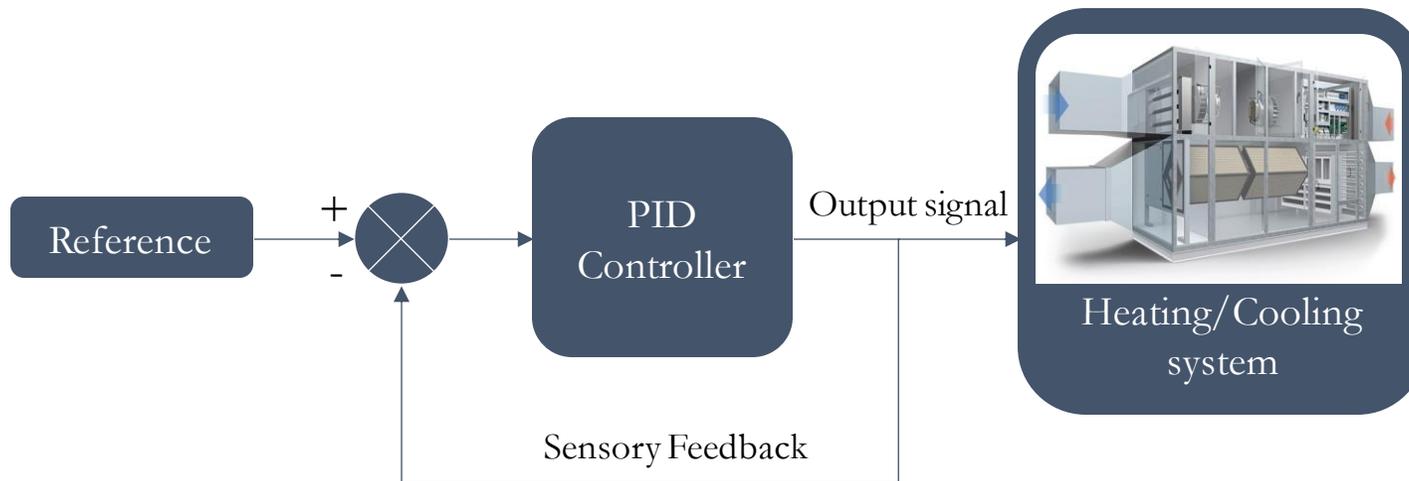
Energy consumption in building

■ Space heating ■ Space cooling ■ Lighting ■ Other



How are building systems controlled today?

A rule-based control system that operates on a fixed set of rules.



Unable to adjust to complex, dynamic environments inside building, for-example: changing thermal demand, price fluctuation, dynamic occupancy.

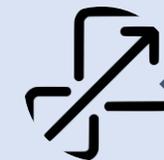
Why RL-based Control?



Adaptive control



Multi-objective Optimization

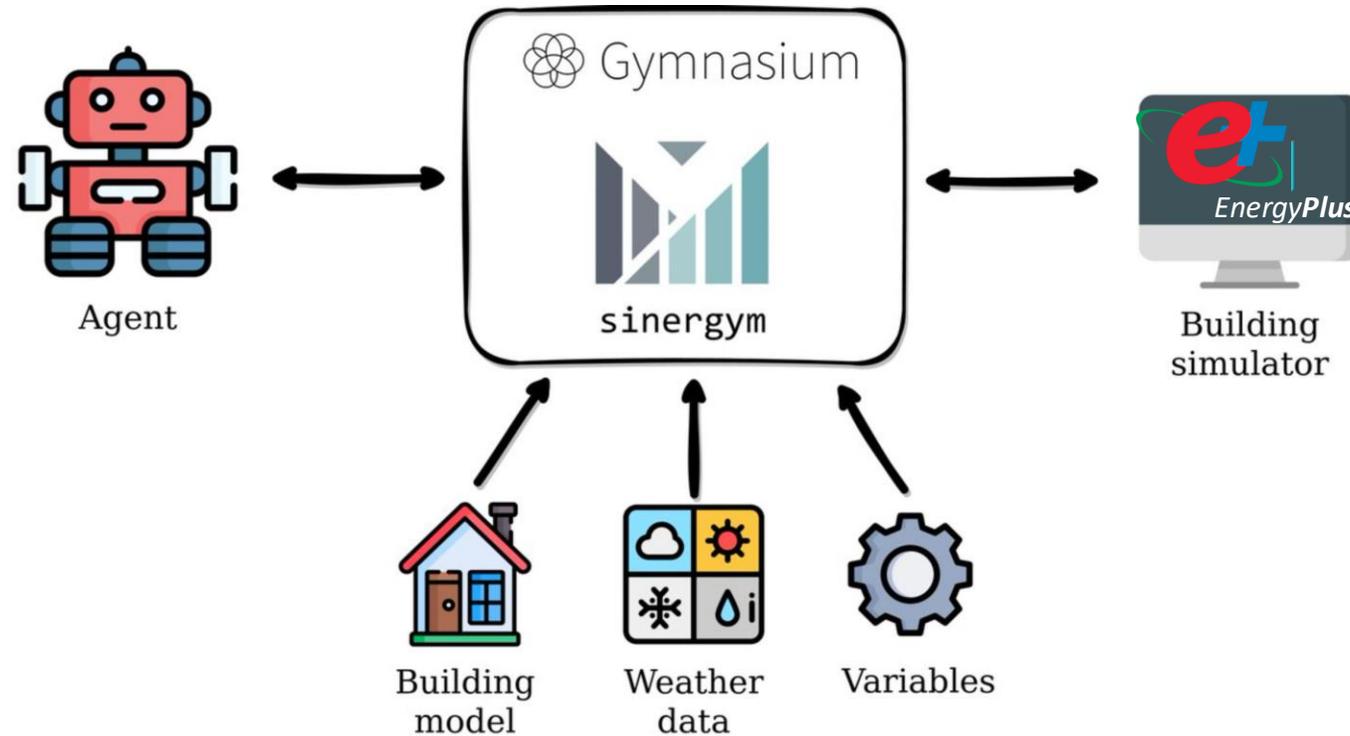


Scalability



Fault detection

Reinforcement Learning control framework



Jiménez-Raboso, J., Campoy-Nieves, A., Manjavacas-Lucas, A., Gómez-Romero, J., & Molina-Solana, M. (2021). Sinergym: A Building Simulation and Control Framework for Training Reinforcement Learning Agents. In *Proceedings of the 8th ACM International Conference on Systems for Energy-Efficient Buildings, Cities, and Transportation* (pp. 319–323). Association for Computing Machinery.

Environment

Input Parameters

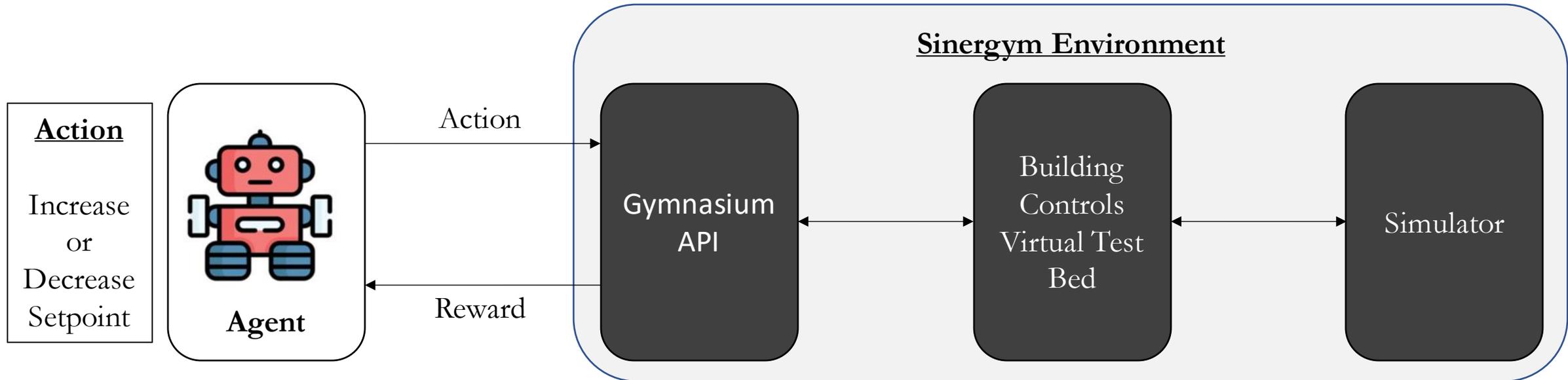
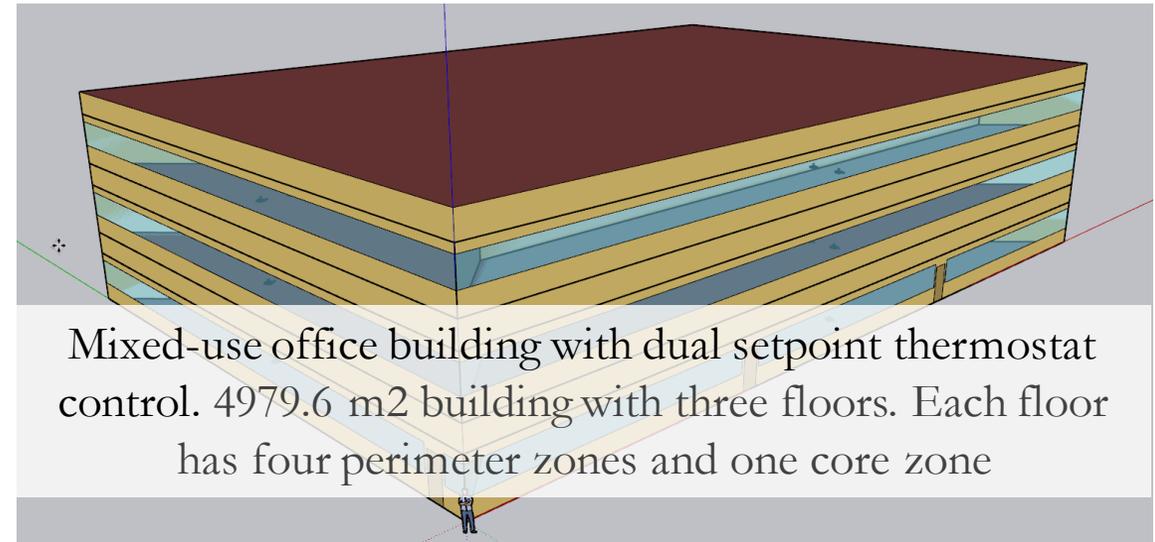
Thermostat heating setpoint,

Thermostat cooling setpoint

Zone Air Temperature (for each zone)

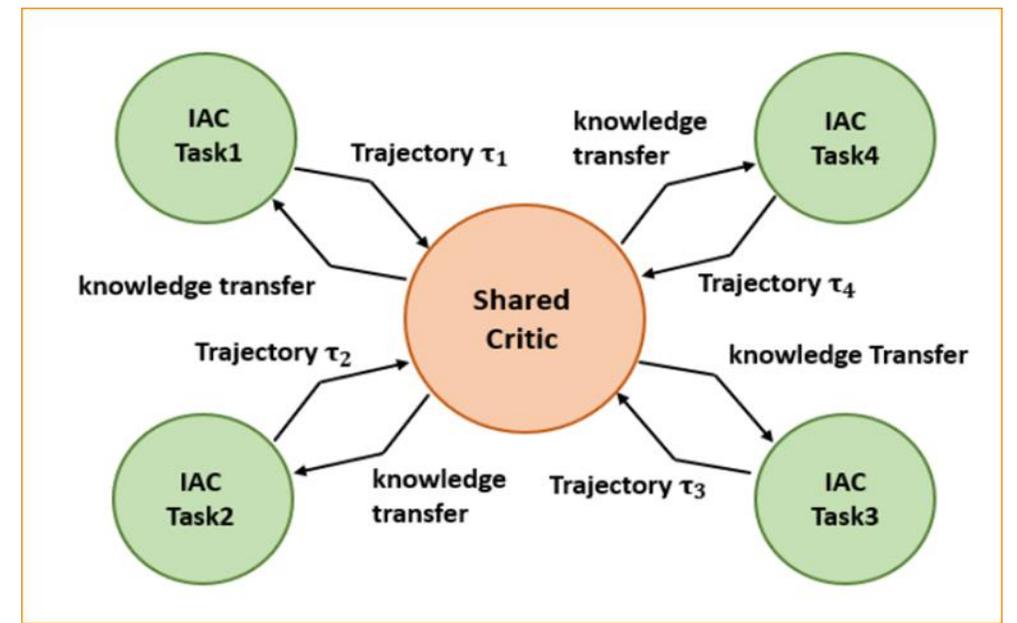
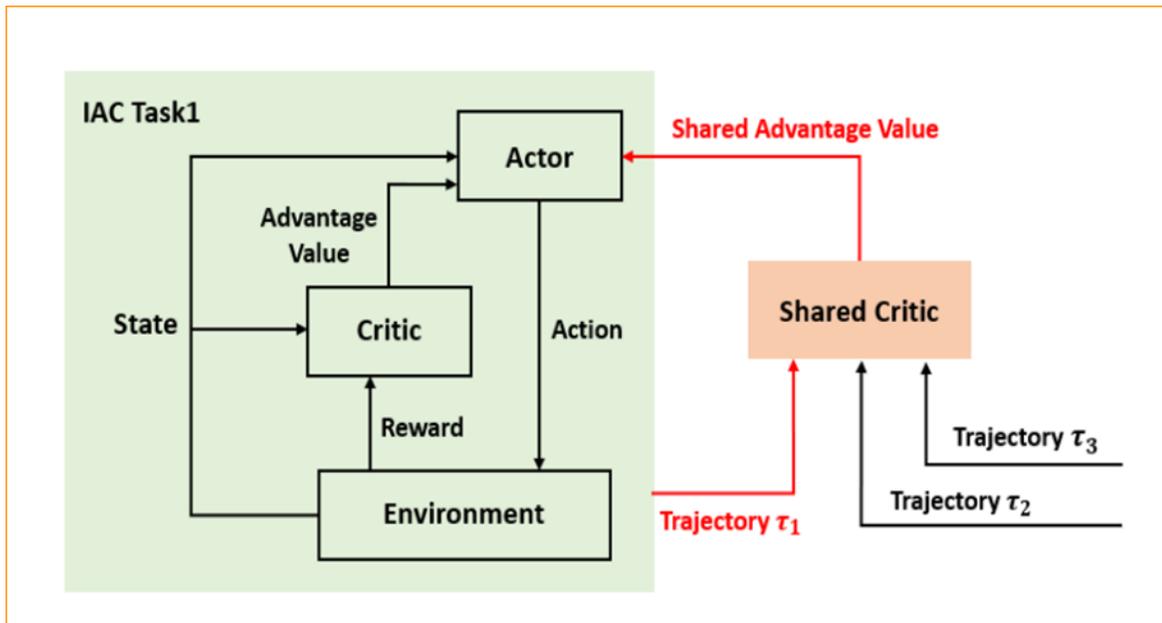
Reward Functions as a trade-off between Energy and Comfort

$$R_t = -w\lambda_E E_t - (1 - \omega)\lambda_T (|T_t - T_{up}| + |T_t - T_{low}|)$$



Implementation

- Algorithm – PPO (baseline) & MultiTaskPPO



Algorithm

Algorithm 1 Multi-task actor-critic with a shared critic

Input: State s ; Reward r

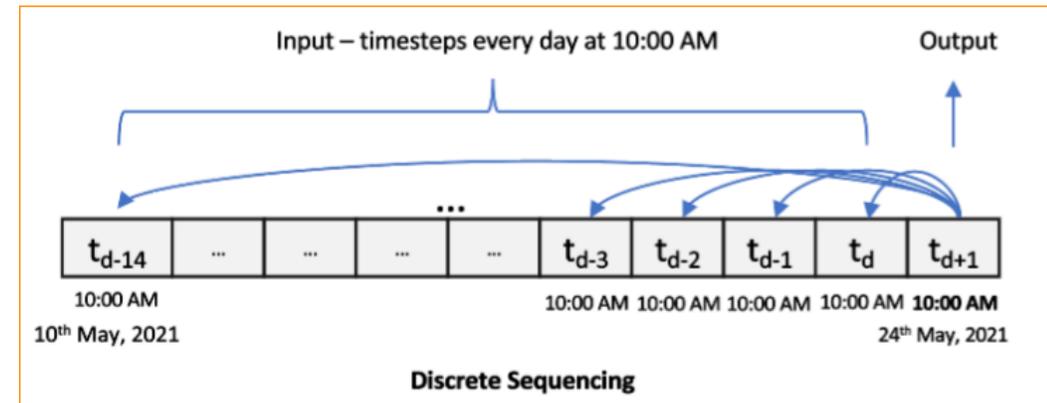
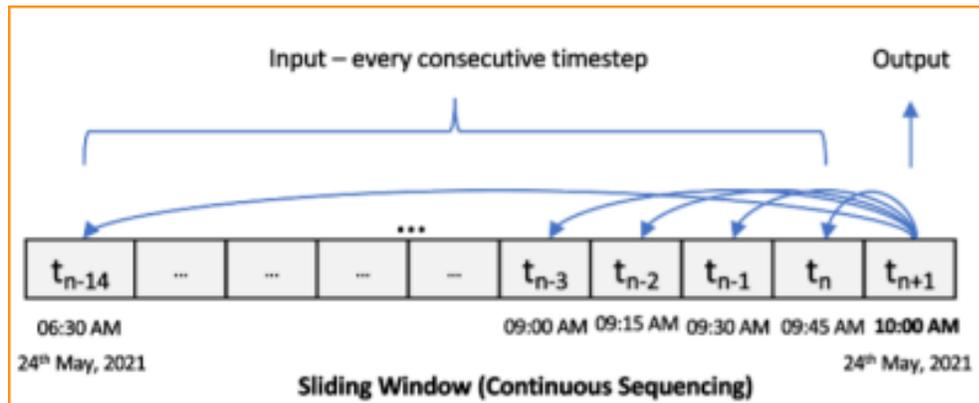
Parameter: Task number $i = 1, 2, \dots, M$; Number of Episode E ; Maximum steps of task per episode T ; Transfer weight α ;

Output: Critic $\eta_{i,E}$ actor $\theta_{i,E}$;

- 1: Randomly initialise actor π_{θ_i} and critic network $V_{\eta_{ni}}$ for task i , $i = 1, 2, \dots, m$;
- 2: Initialize episode counter e , $e = 0$
- 3: **while** $e_i < E_i$ **do**
- 4: **for** each task i **do**
- 5: Set step counter $t = 0$
- 6: **while** $t < T_i$ and not terminal state **do**
- 7: Select action $a_{i,t} \sim \pi_{\theta_i}$.
- 8: Execute $a_{i,t}$ and state $r_{i,t+1}$ and $s_{i,t+1}$
- 9: Store tuple $(s_{i,t}, a_{i,t}, r_{i,t+1}, s_{i,t+1})$
- 10: Update step counter: $t \leftarrow t + 1$
- 11: **end while**
- 12: **for** each sample **do**
- 13: Compute advantage value using Eq.10
- 14: Compute shared advantage value using Eq.8
- 15: **end for**
- 16: Compute critic gradient using Eq.12
- 17: Compute actor gradient using Eq.11
- 18: Compute shared critic gradient using Eq.6
- 19: **end for**
- 20: Update episode counter $e \leftarrow e + 1$
- 21: **end while**
- 22: **return** Critic and actor weights: $\theta_{i,E}, \eta_{i,E}$;

Models

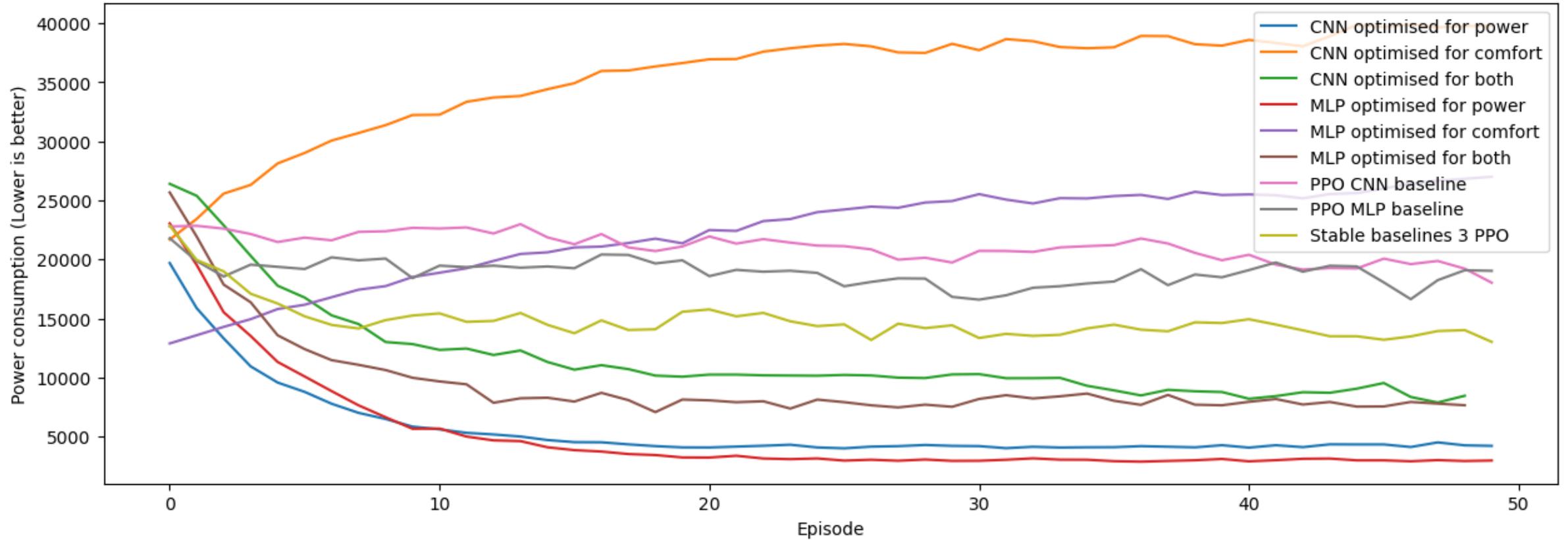
- MLP
- Time-discrete CNN (1D)



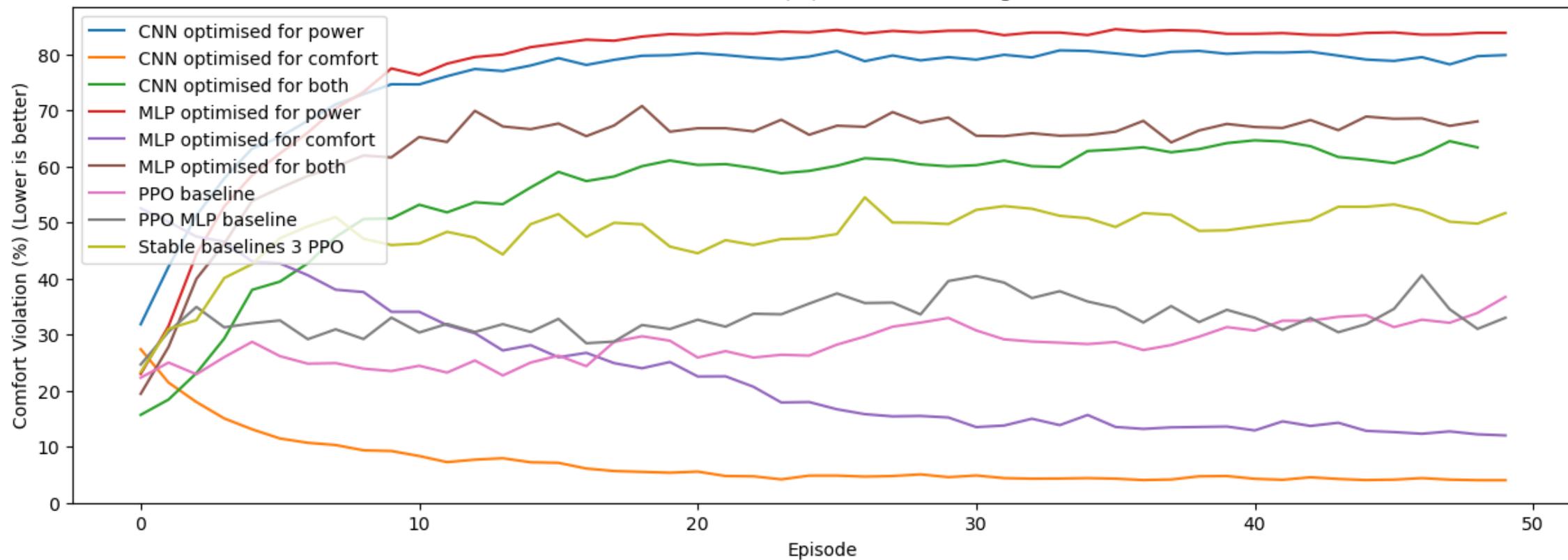
Nisha Menon, Shantanu Saboo, Tanmay Ambadkar, & Umesh Uppili (2022). Discrete Sequencing for Demand Forecasting: A novel data sampling technique for time series forecasting. In *3rd International Conference on Intelligent Data Science Technologies and Applications, IDSTA 2022, San Antonio, TX, USA, September 5-7, 2022* (pp. 61–67). IEEE.

Results

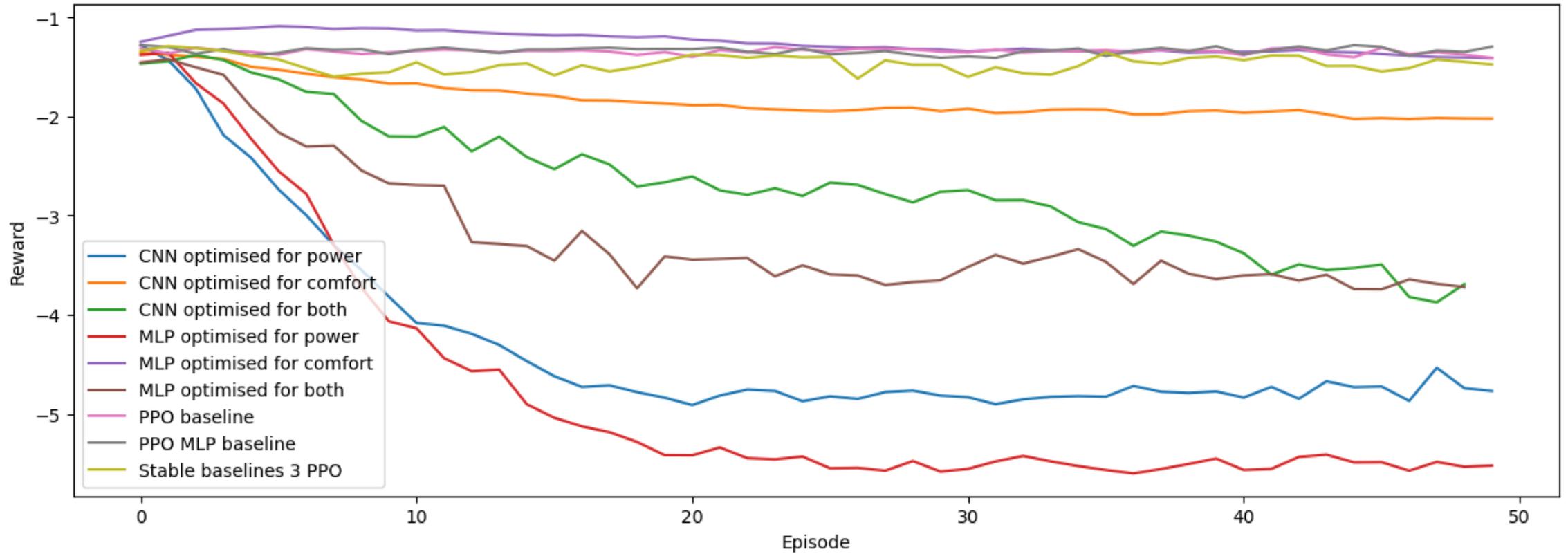
Power consumption of different RL agents



Comfort Violation (%) of different RL agents



Reward of different RL agents



Thank You!