

Artificial Intelligence for HVAC Systems

V. Ziebart, ACSS IAMP ZHAW

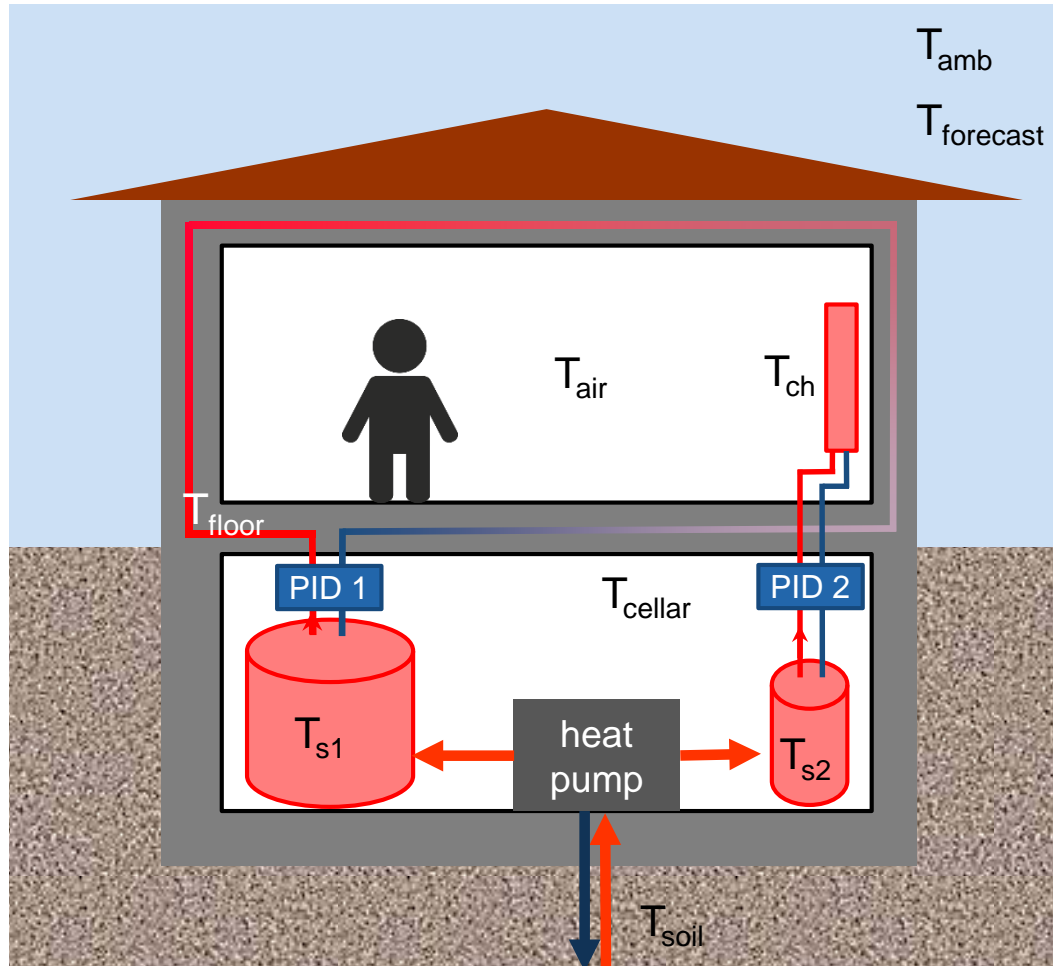
Motivation

- More than 80% of the energy consumed in Swiss households is used for room heating & cooling and domestic hot water.*
- The increasing combined use of various energy conversion and storage technologies (PT, solar thermal collectors, heat pumps, combustion, batteries, hot water storage, ice storage systems) requires intelligent and optimal control systems.
- Reinforcement Learning (RL) showed promising performance in different fields (AlphaZero, computer games).
- RL is a data-driven approach.
- Can RL be used as control method for HVAC**-Systems?
(optimality, learning behavior, robustness,...

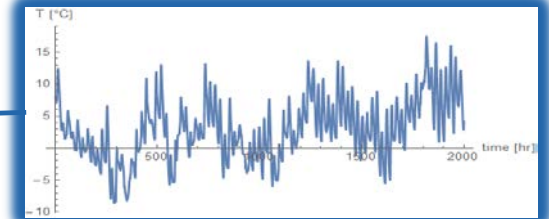
* Energieverbrauch in der Schweiz und weltweit, EnergieSchweiz, Bundesamt für Energie BFE
Dienst Aus-und Weiterbildung, Juli 2015

** HVAC: heating, ventilation, air conditioning

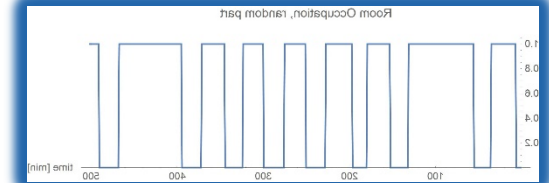
Model of Building and Heating System



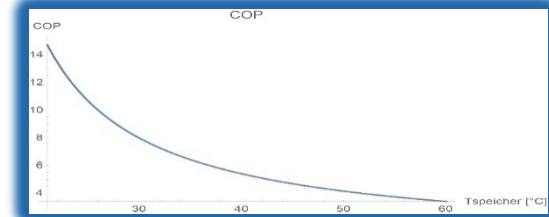
Weather data of 11 Swiss cities



Occupation pattern



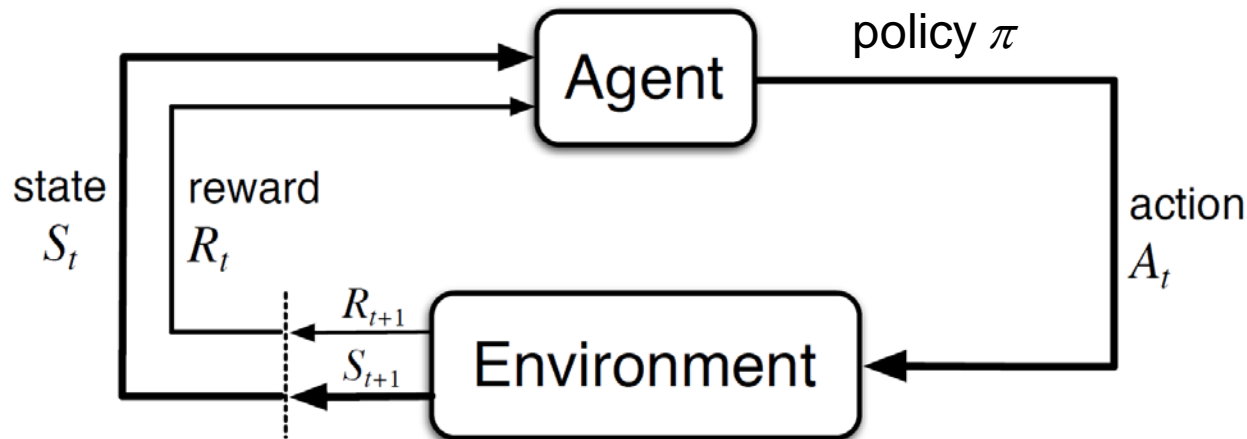
COP of heat pump



Energy price

night rate	0.5
date rate	1

Reinforcement Learning Control



value function: $q_{\pi}(s, a) = \mathbb{E}_{\pi} [G_t | S_t = s, A_t = a]$

return: $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$

discount rate: $\gamma \in [0, 1]$

Q-Learning and SARSA

Q-learning:

$$q(s_t, a_t) \leftarrow q(s_t, a_t) + \alpha \left(\overbrace{R_t + \gamma \max_a q(s_{t+1}, a)}^{\text{new estimate of return G}} - q(s_t, a_t) \right)$$

↑
learning rate

SARSA:

$$q(s_t, a_t) \leftarrow q(s_t, a_t) + \alpha \left(R_t + \underbrace{\gamma q(s_{t+1}, a_{t+1})}_{\text{difference of «large» numbers!}} - q(s_t, a_t) \right)$$

These *interacting, interdependent, iterative* jobs must be done:

1. learn q for actions performed in state visited
2. approximate a generalized q -function on state and action space (SxA) (typically with gradient descent)
3. improve control performance by e.g. (ε -)greedy policy

→ **potentially unstable**

Improvement of Stability

- Ad 1: More efficient learning of q through n -step SARSA

$$G_{t,n} = \underbrace{\sum_{k=0}^{n-1} \gamma^k R_{t+k}}_{\text{truncated sum of } n \text{ rewards}} + \gamma^n q(s_{t+n}, a_{t+n})$$

instead of $G_{t,1} = R_t + \gamma q(s_{t+1}, a_{t+1})$.

Average fraction return:
$$\frac{\sum_{k=0}^{n-1} \gamma^k \bar{R}}{\sum_{k=0}^{\infty} \gamma^k \bar{R}} = \frac{1 - \gamma^n}{1 - \gamma} = 1 - \gamma^n$$

- Ad 2: Least-Squares fit to $\{s_t, a_t, G_{t,n}\}$ with polynomials and trigonometric functions to find $q(s, a)$ on $S \times A$ instead of iterative gradient descent methods.

State, Action, Reward

- **State** variables for RLC (continuous and discrete)
 - **6 temperatures**, \mathbb{R}^6 (T_{amb} , T_{air} , T_{floor} , T_{storage1} , T_{storage2} , T_{forecast})
 - **time**, real interval $[0,24[$
 - **room occupancy**, boolean
- **Action** variables (discrete, 12 combinations)
 - **heat pump off/loading storage 1 or 2**, $\in \{0,1,2\}$
 - **PID floor heating on/off**, boolean
 - **PID convection heater on/off**, boolean
- **Reward** (≤ 0):
 - **energy costs** (negative, night and day rate)
 - **temperature deviation** from setpoint, (negative, proportional to ΔT , only if house is occupied)

Approximation of State-Action Value Function $q(s,a)$

- Partitioning of state-action space in continuous (c) and discrete (d) subspaces:

$$S \times A = (S_c \times S_d) \times (A_c \times A_d) = \underbrace{S_c \times A_c}_{\text{continuous}} \times \underbrace{S_d \times A_d}_{\text{discrete}}$$

- For each set of discrete variables in $S_d \times A_d$ (24 configurations) $q(s,a)$ is approximated in the continuous subspace $S_c \times A_c$ by a linear combination of polynomials (temperatures) and trigonometric functions (time):

$$q_k(s, a) = \sum_{j,l} \tilde{c}_{kjl} \prod_{i=1}^6 T_i^{n(k,i,j)} \text{trig}_l(t) \quad k = 1, \dots, 24$$

rewritten as flattened vector product

$$q_k(s, a) = \sum_j c_{kj} p_j(T_1, \dots, T_6, t) = \mathbf{c}_k \cdot \mathbf{p} \quad k = 1, \dots, 24$$

typically: $j = 252 \rightarrow k \cdot j = 6'048$ coefficients c_{kj}

Analytical Solution

older experience is discounted by γ_Q

$$\mathbf{c}_k = \arg \min \sum_{i=1}^t \gamma_Q^{t-i} \left\{ \left(\sum_{j=0}^{n-1} \gamma^j R_{s(i,k)+j} + \gamma^n q(s_{s(i,k)+n}, a_{s(i,k)+n}) \right) - \mathbf{c}_k \mathbf{p}_{s(i,k)} \right\}^2, \quad k = 1, \dots, 24$$

$s(i,k)$ is the index when $(s_d, a_d)_{i-1} = (s_d, a_d)_k$ the i -th time

rearranging yields:

$$\mathbf{Q}_k^{(i)} = \gamma_Q \mathbf{Q}_k^{(i-1)} + \left\{ \mathbf{p}_{s(i,k), r} \mathbf{p}_{s(i,k), l} \right\}_{rl}$$

$$\mathbf{b}_k^{(i)} = \gamma_Q \mathbf{b}_k^{(i-1)} - 2 \mathbf{p}_{s(i,k)} \left(\sum_{j=0}^{n-1} \gamma^j R_{s(i,k)+j} + \gamma^n q(s_{s(i,k)+n}, a_{s(i,k)+n}) \right)$$

$$\Rightarrow \mathbf{c}_k = 0.5 \cdot (\mathbf{Q}_k + \varepsilon_k \mathbf{I})^{-1} \mathbf{b}_k$$

for numerical reasons

Decision Making, Reward Normalization

- probability of choosing a_j :

$$\pi(a_i | s) = \frac{e^{q(s, a_i)/\tau}}{\sum_{j=1}^n e^{q(s, a_j)/\tau}}$$

$$\tau \rightarrow \infty : \pi(a_i | s) = \pi(a_j | s)$$

$$\tau \rightarrow 0 : \pi(\arg \max q(s, a_i) | s) = 1$$

$$q(s, a_i) \leq 0!$$

- normalization of reward per year (plots only!) by difference of average outside temperature and room set point temperature

$$R_{year, norm} = \frac{R_{year}}{\Delta T + 0.005 \Delta T^2} \quad \text{with } \Delta T = T_{sp} - \bar{T}_{amb}$$

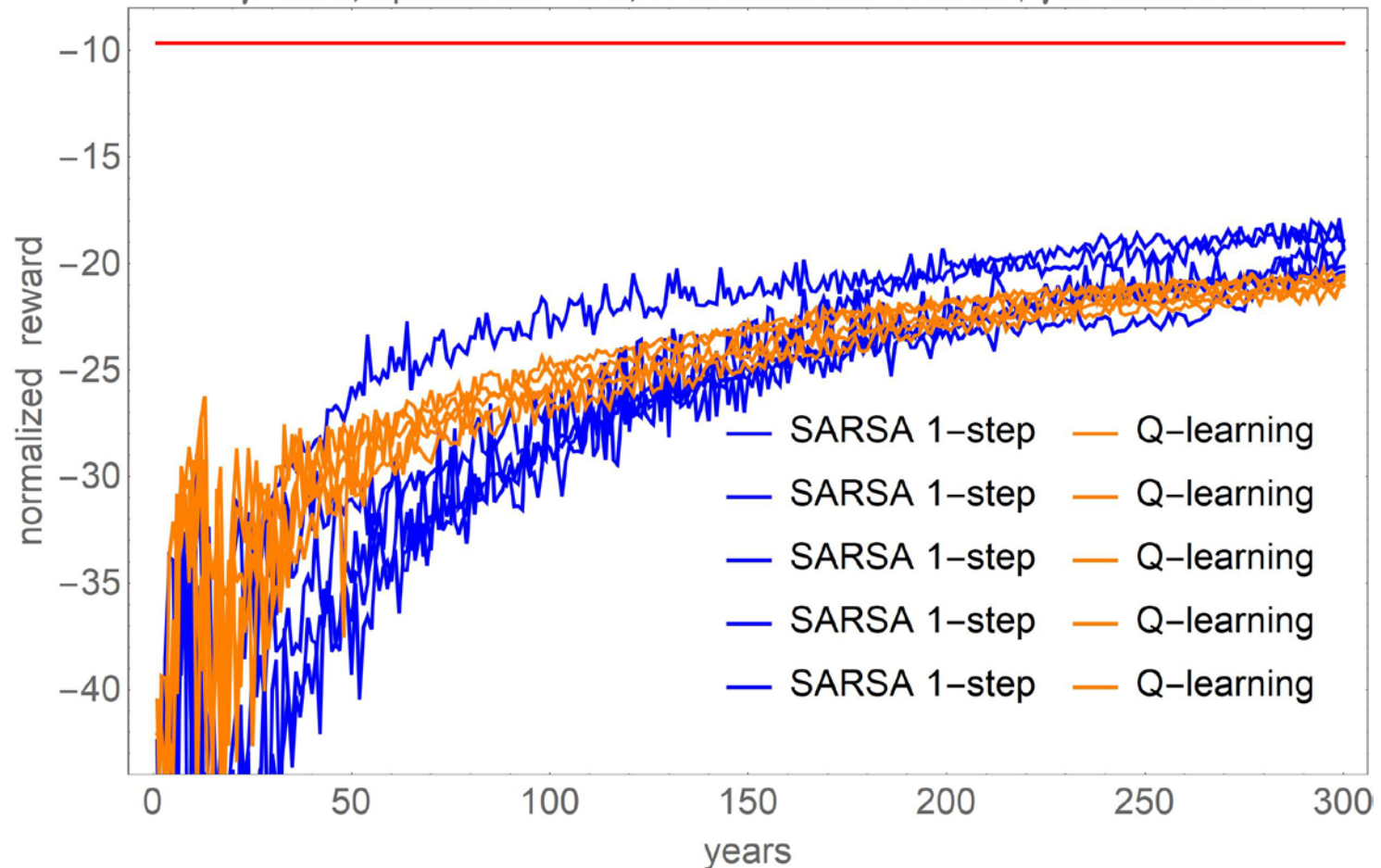
Hardware & Software

- 2 Servers with 2 CPUs Intel Xeon Platinum 8164 each
- Each CPU with 26 cores
- 768 GB RAM
- Code written in Mathematica
- Computation time for 1 year simulation: 1min

SARSA vs Q-Learning

Comparison of Q-learning with 1-step SARSA

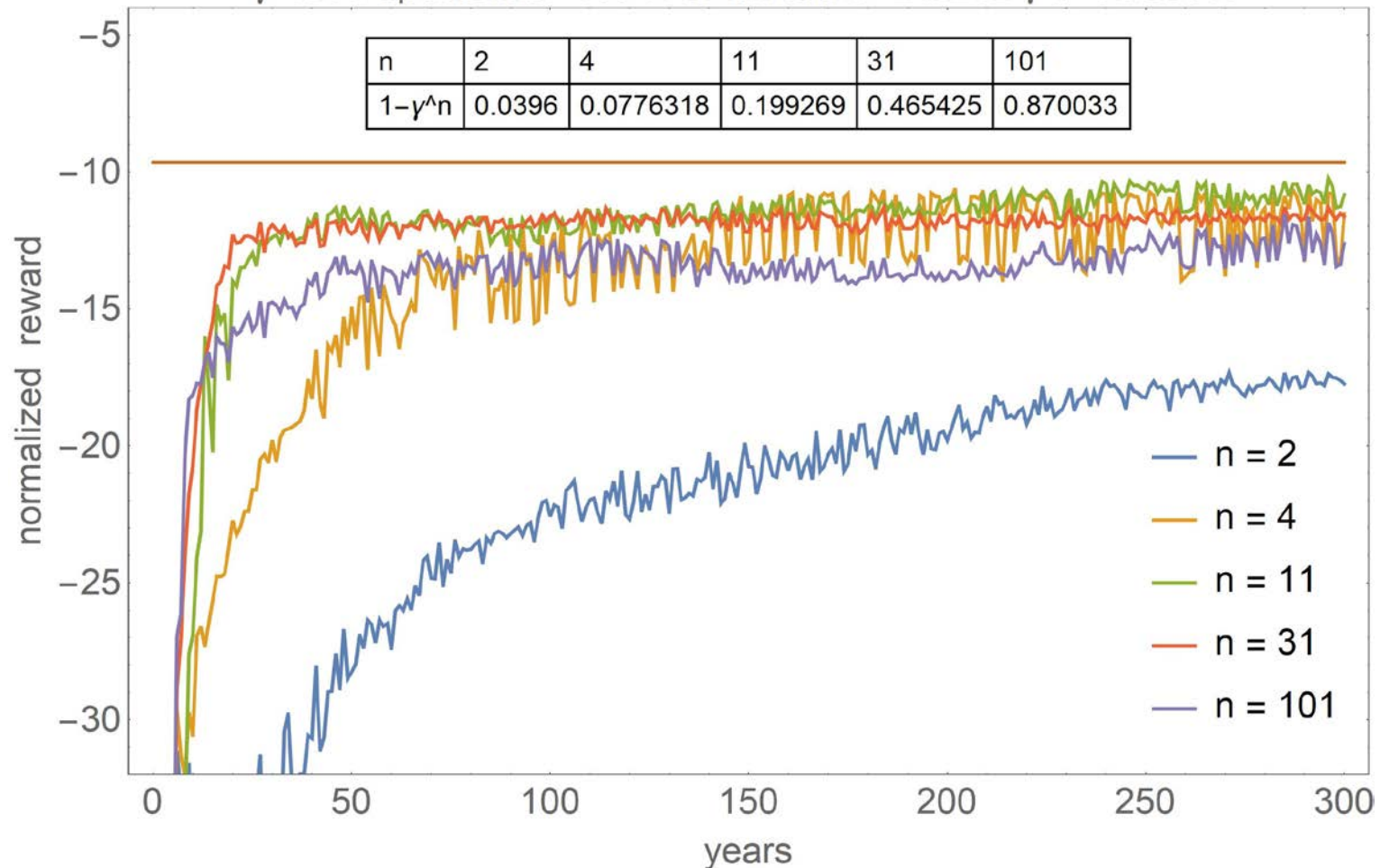
$\gamma=0.98$, updatrate=400, # of functions=24x252, $\gamma_Q=0.999987$



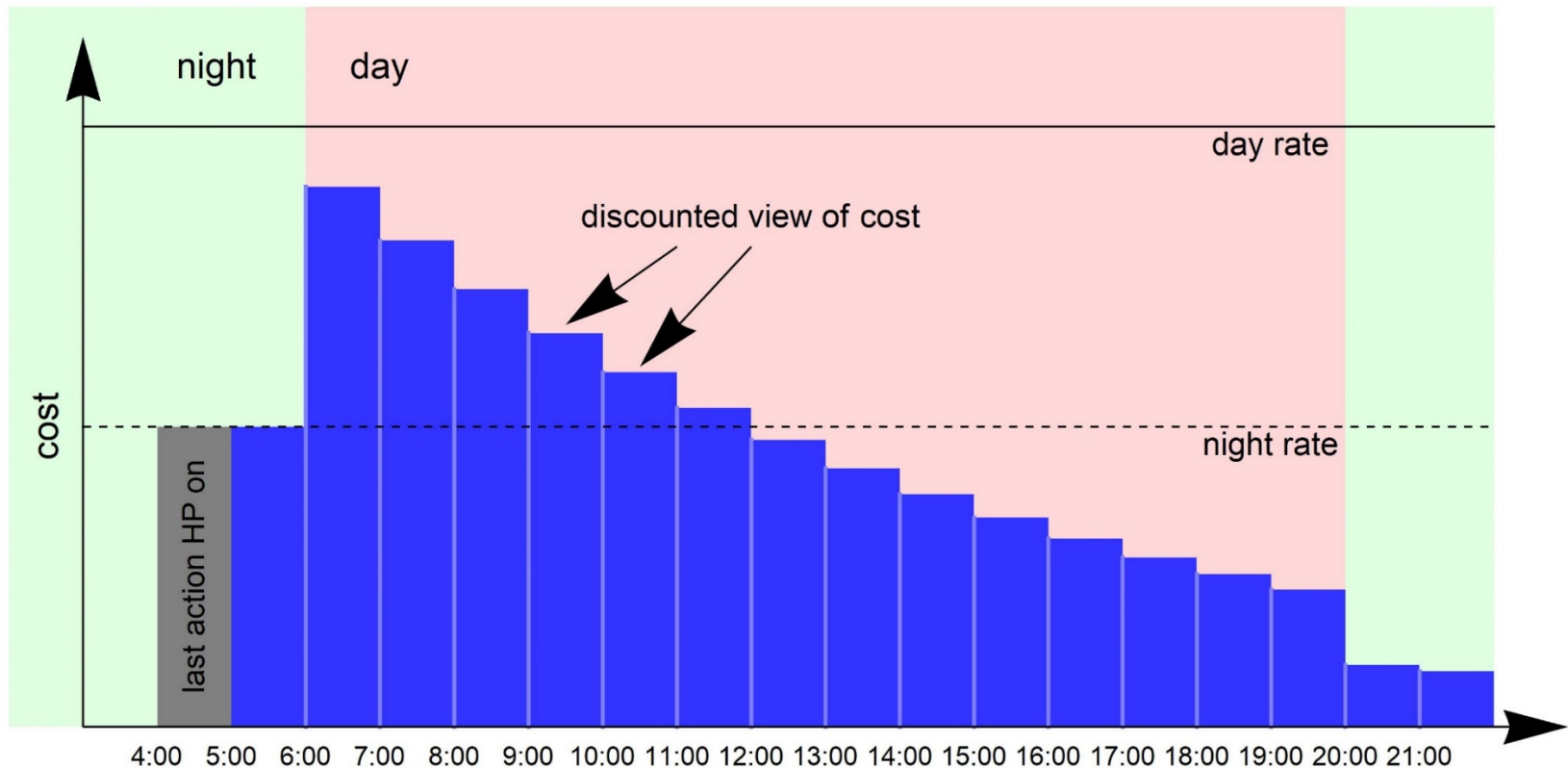
Effect of n in n -step SARSA

Learning Behaviour of n -step SARSA

$\gamma=0.98$ updatrate=400 # of functions=24x252 $\gamma Q=0.999987$



Influence of Discount Factor γ



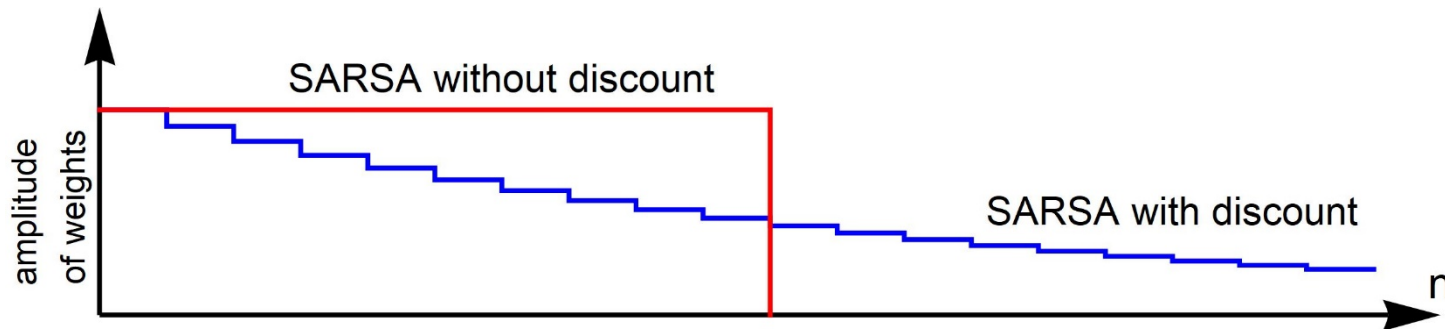
Actual time 5:00. One more heating hour before 20:00 is needed. What hour seems to be the best choice? (If $\gamma > 0.952$: heating at 5:00, otherwise heating at 19:00)

→ especially too low γ 's lead to suboptimal decisions and unrealistic q -values

n -step SARSA without Discount

- The discount of future rewards disturbs the optimal scheduling of actions with fixed and known costs.
- Alternative approach: n -step SARSA without discount

$$G_{t,n} = \sum_{k=0}^{n-1} \cancel{\gamma^k} R_{t+k} + \cancel{\gamma^n} \cancel{q(s_{t+n}, a_{t+n})}$$

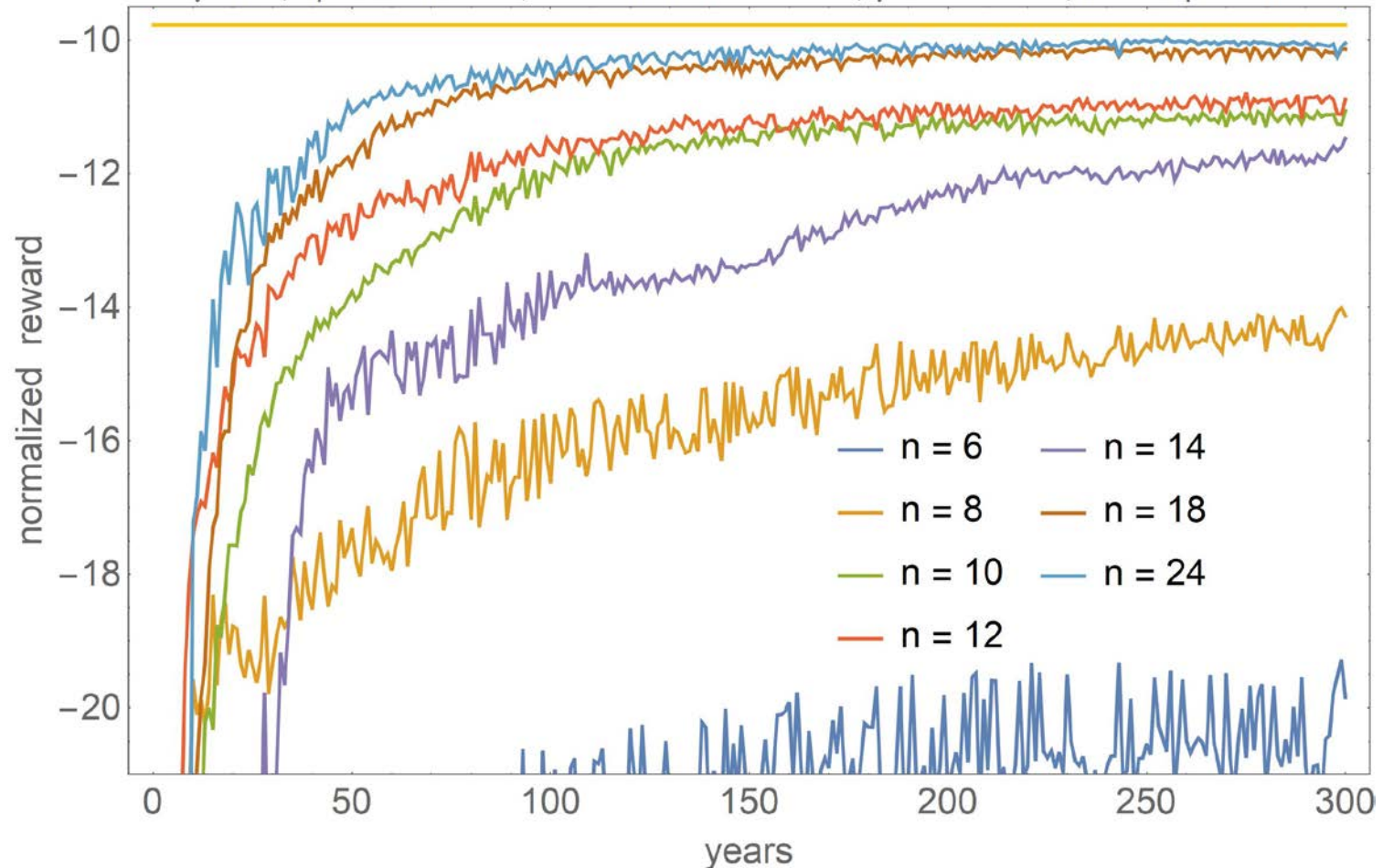


- The time horizon is defined by n instead of γ !

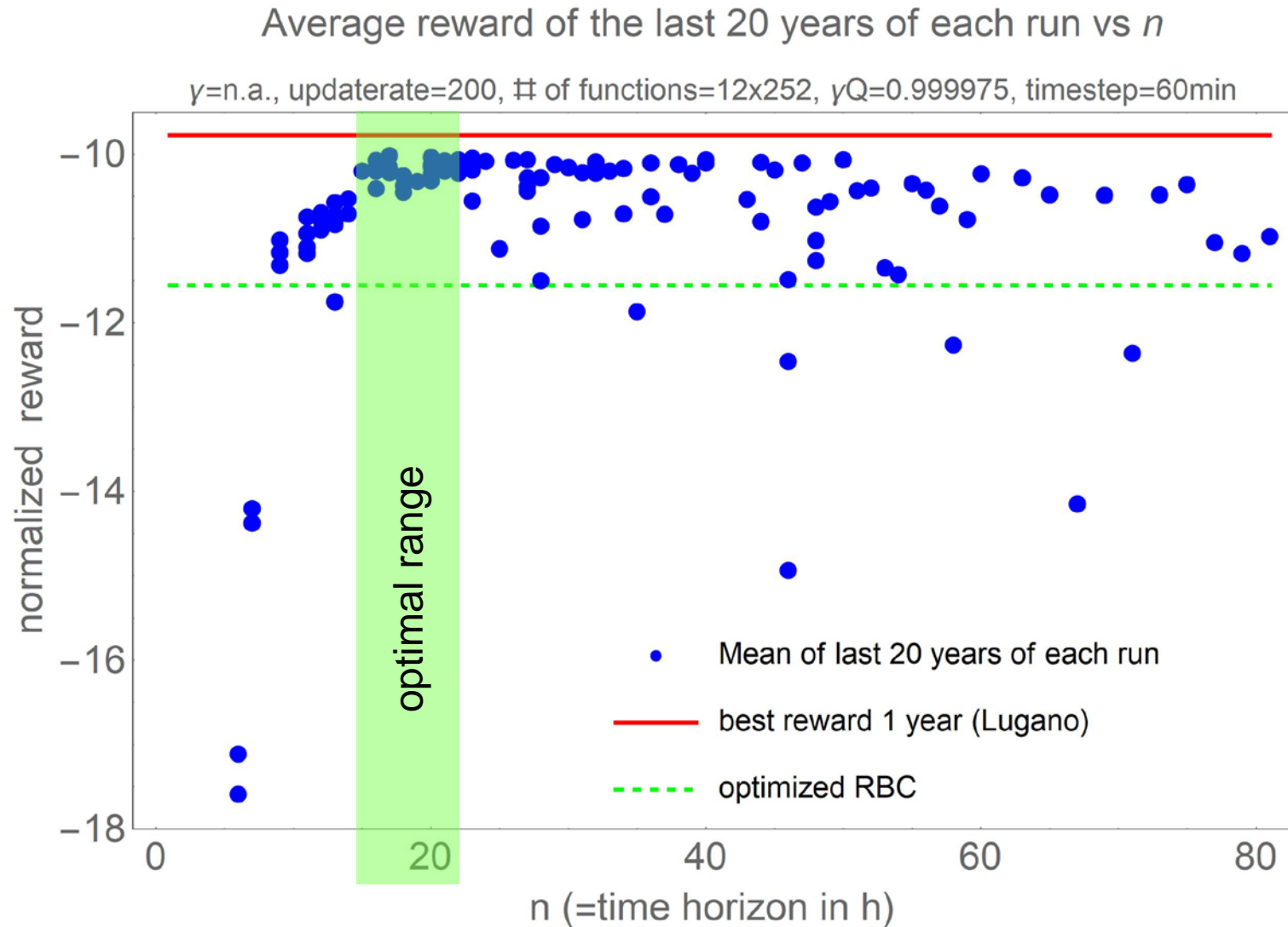
n -step SARSA without Discount

Learning Behaviour of n -step SARSA

$\gamma = \text{n.a.}$, $\text{update rate} = 200$, $\# \text{ of functions} = 12 \times 252$, $\gamma Q = 0.999975$, $\text{timestep} = 60 \text{ min}$



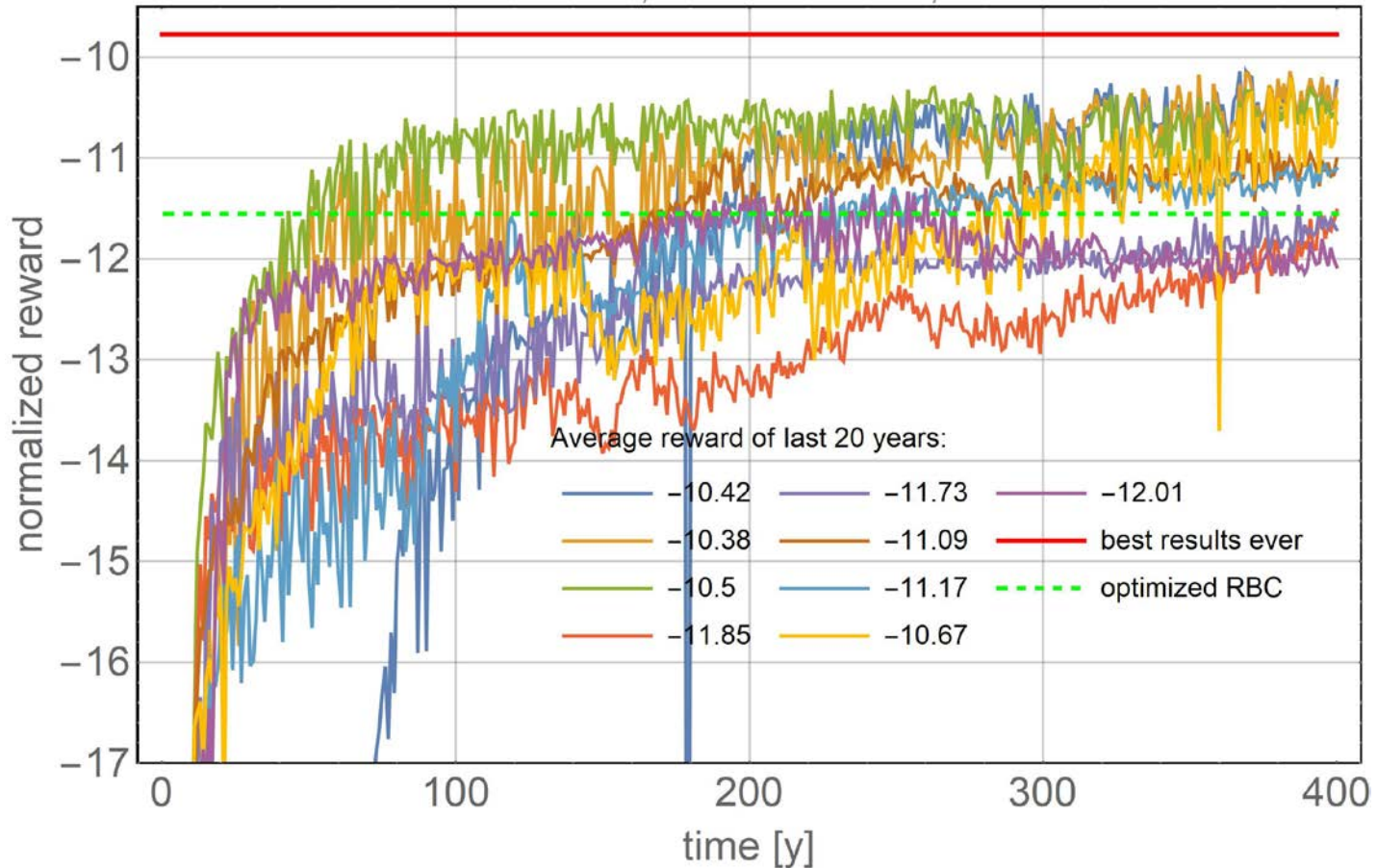
Performance of n -step SARSA without Discount



Parallel «Supporting» Agents, Motivation

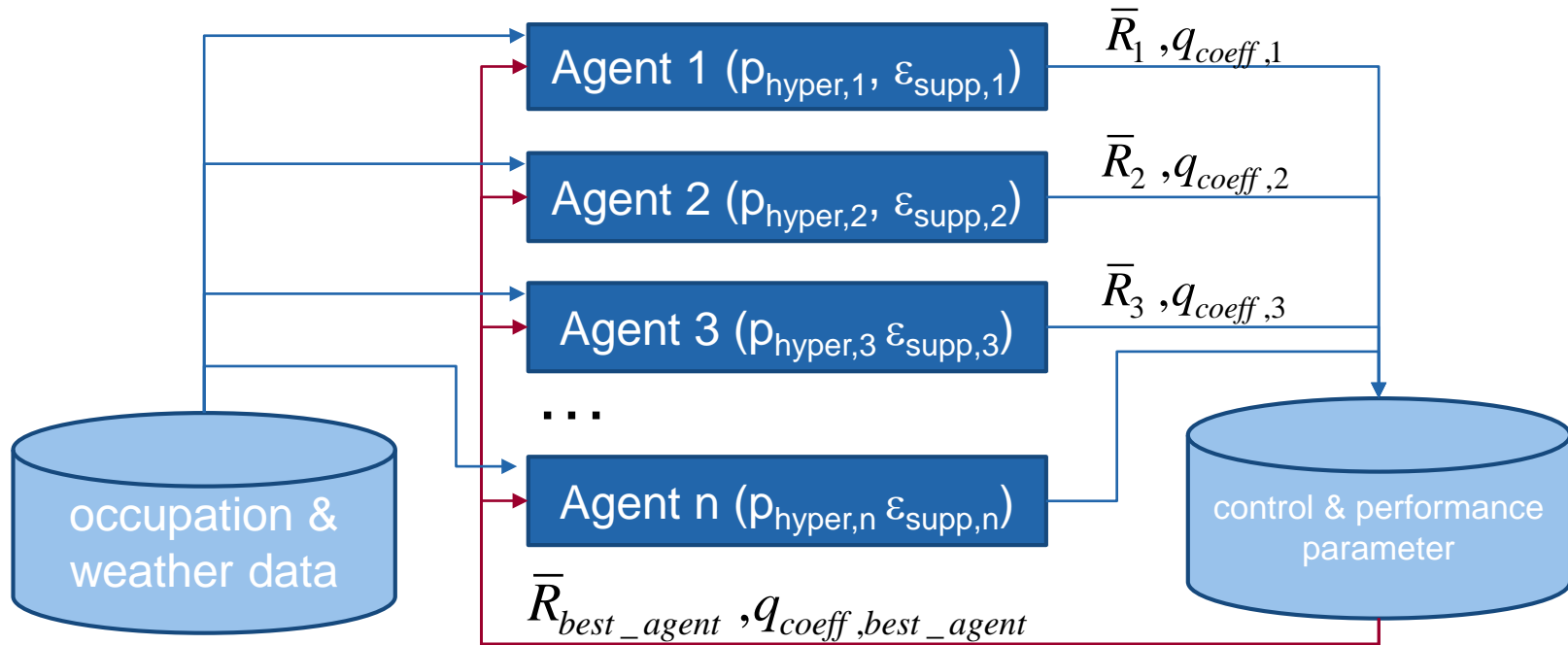
Learning behaviour of 9 agents (different hyperparameters)

no discount, #functions=24*252, n=48



Could larger scatter be exploited by helping worse agents become better and then by chance even be better than best agent?

Parallel «Supporting» Agents, Structure



$$\varepsilon_{supp,i} = \begin{cases} \min \left[c_{supp} \frac{\bar{R}_i}{(\bar{R}_i - \bar{R}_{best_agent})} \left(1 - e^{\#year/\tau_{supp}} \right), 1 \right] & \text{if \# year even} \\ 0 & \text{if \# year odd} \end{cases}$$

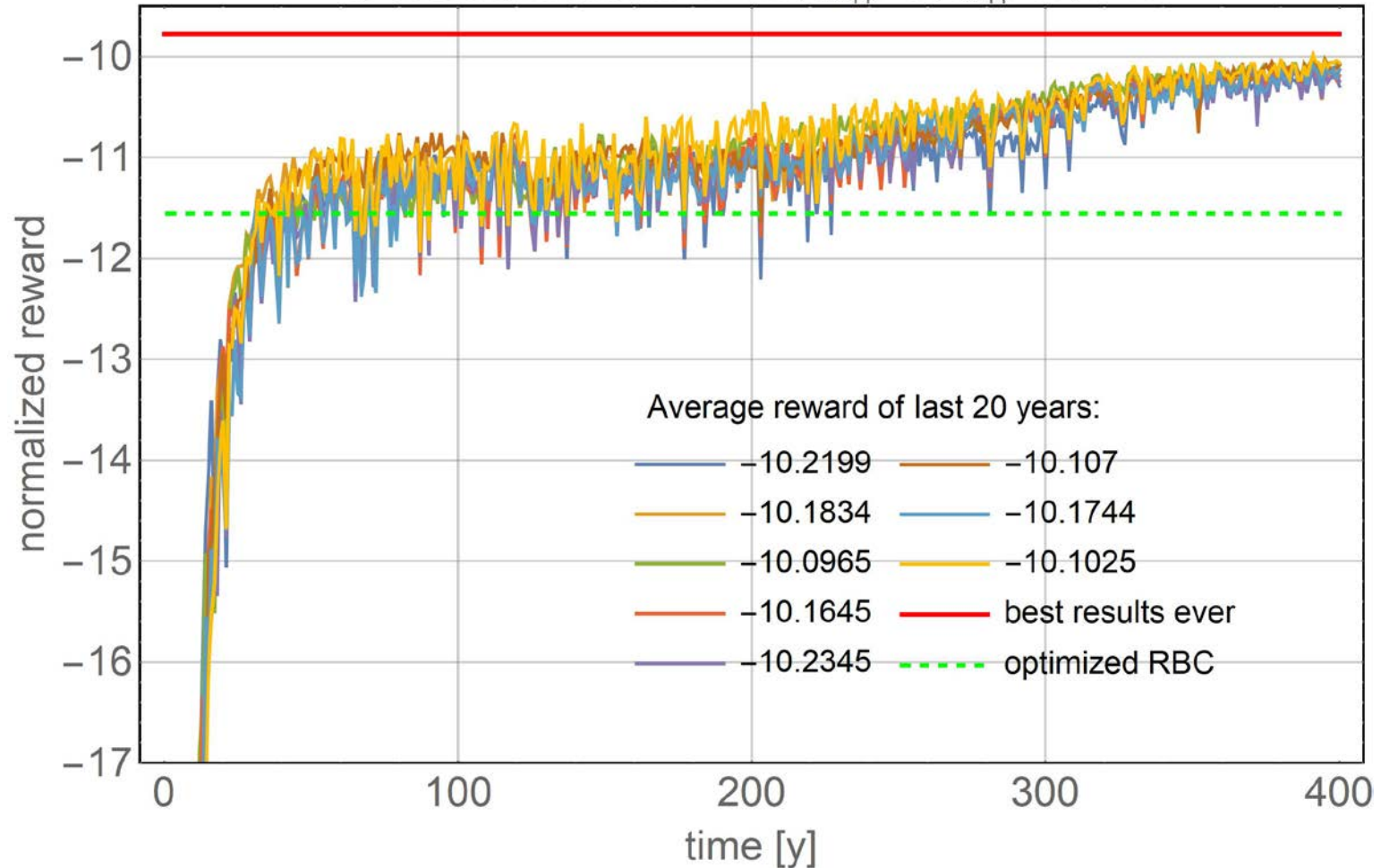
needed to evaluate performance ←

$\varepsilon > \varepsilon_{supp}$: decision based on own parameters

$\varepsilon < \varepsilon_{supp}$: decision based on best agents parameters

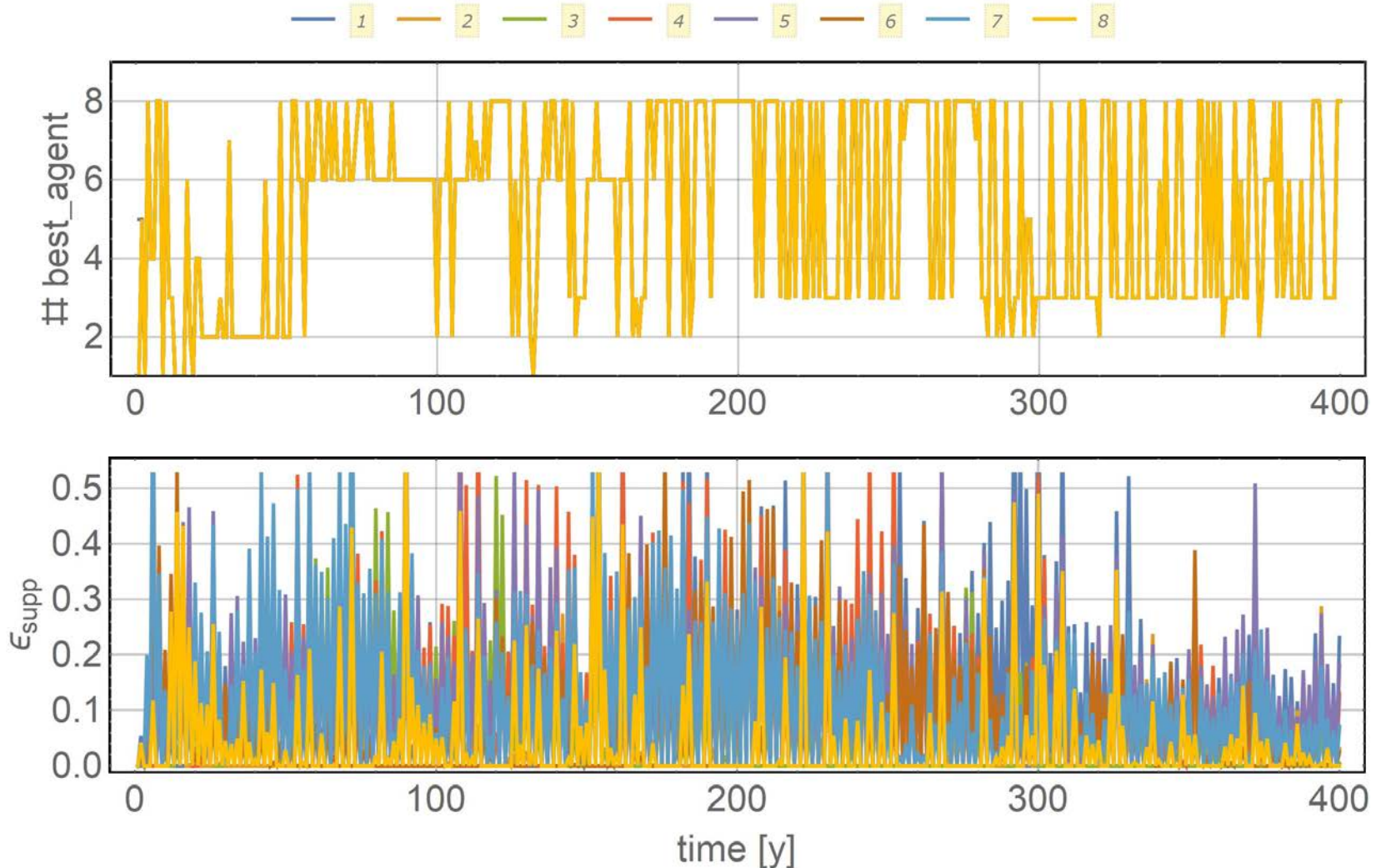
Parallel «Supporting» Agents, Results

Learning behaviour of 8 mutually supporting agents (same hyperp.)
no discount, #functions=24*252, $c_{\text{supp}}=10$, $\tau_{\text{supp}}=30$, $n=48$



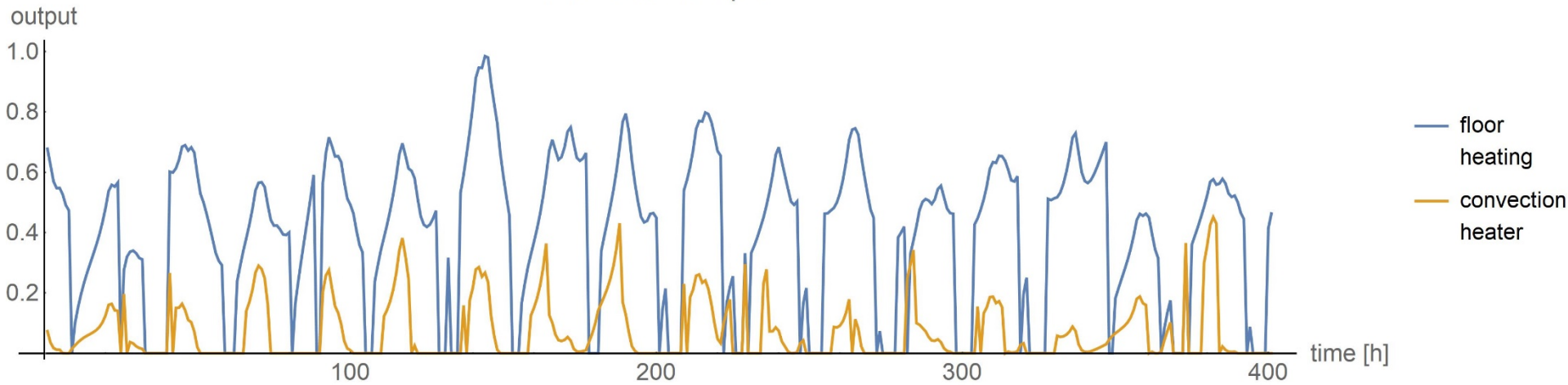
→ all agents show better performance than best agent without support

Parallel «Supporting» Agents, Results

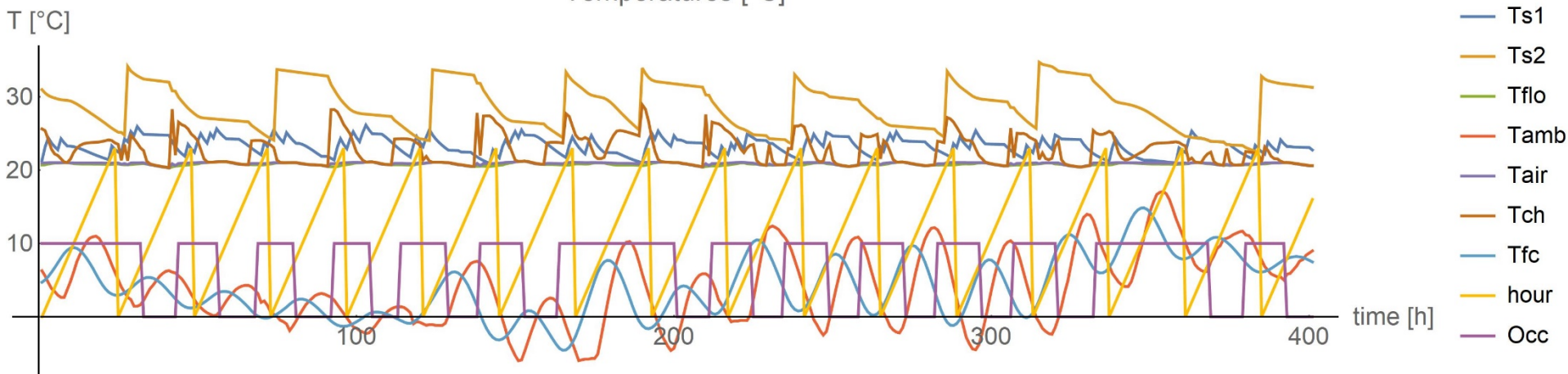


RL Control Example (1)

PID control output

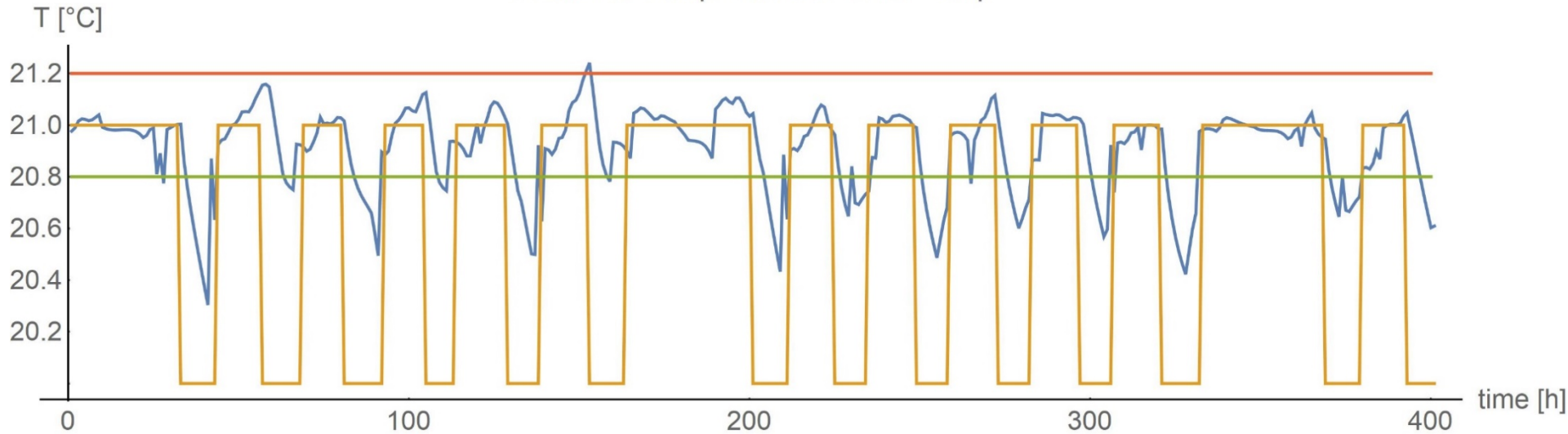


Temperatures [°C]

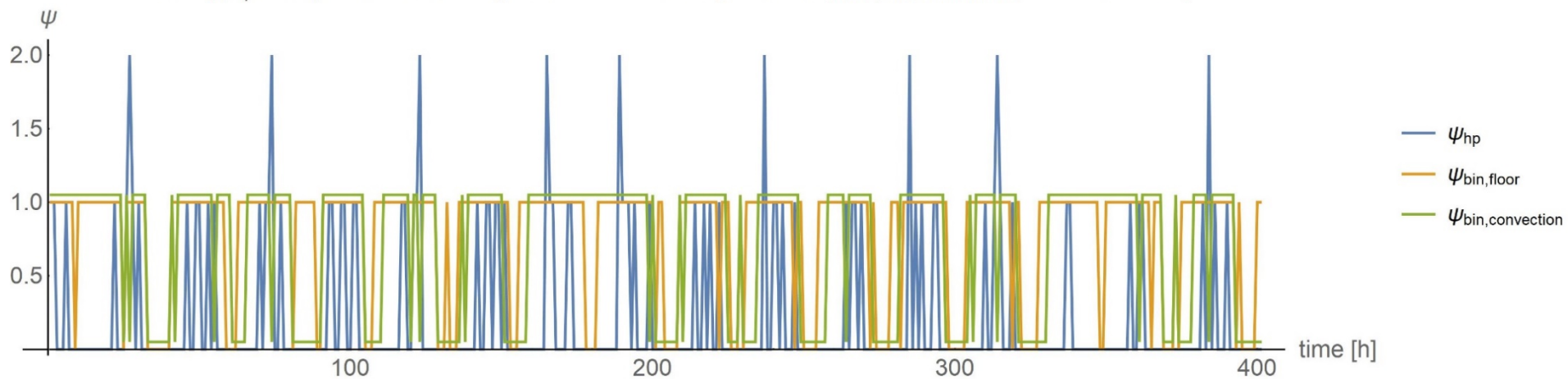


RL Control Example (2)

Room Air Temp. and Set Point Temp

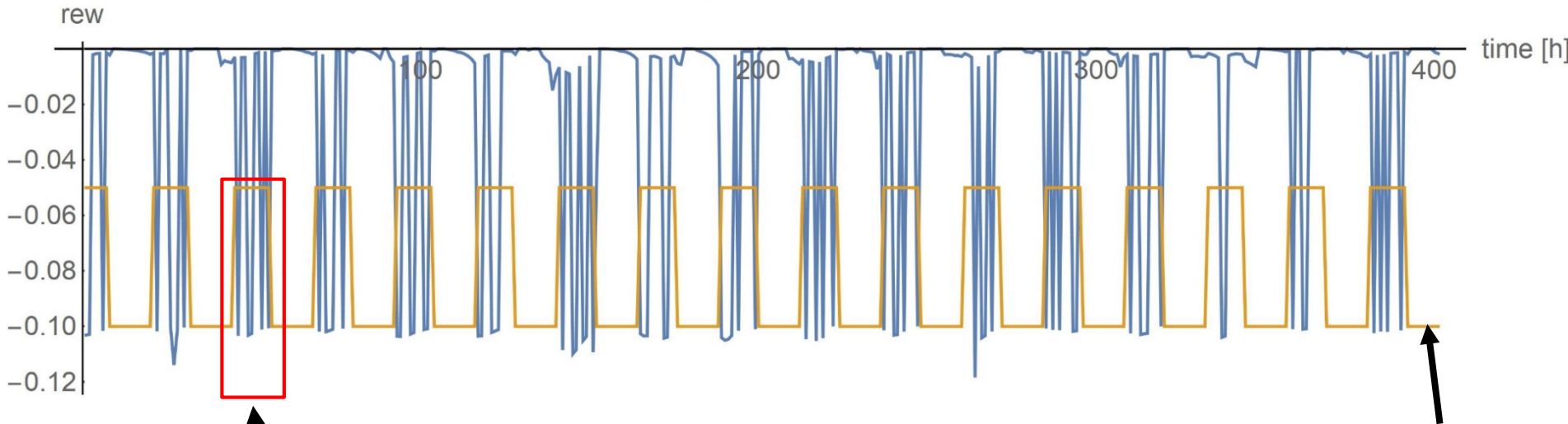


Actions [$\psi_{hp}=0$: hp off, =1: loading tank 1, =2: loading tank 2; $\psi_{floor/convection\ heater}=0$: off, =1 on]



RL Control Example (3)

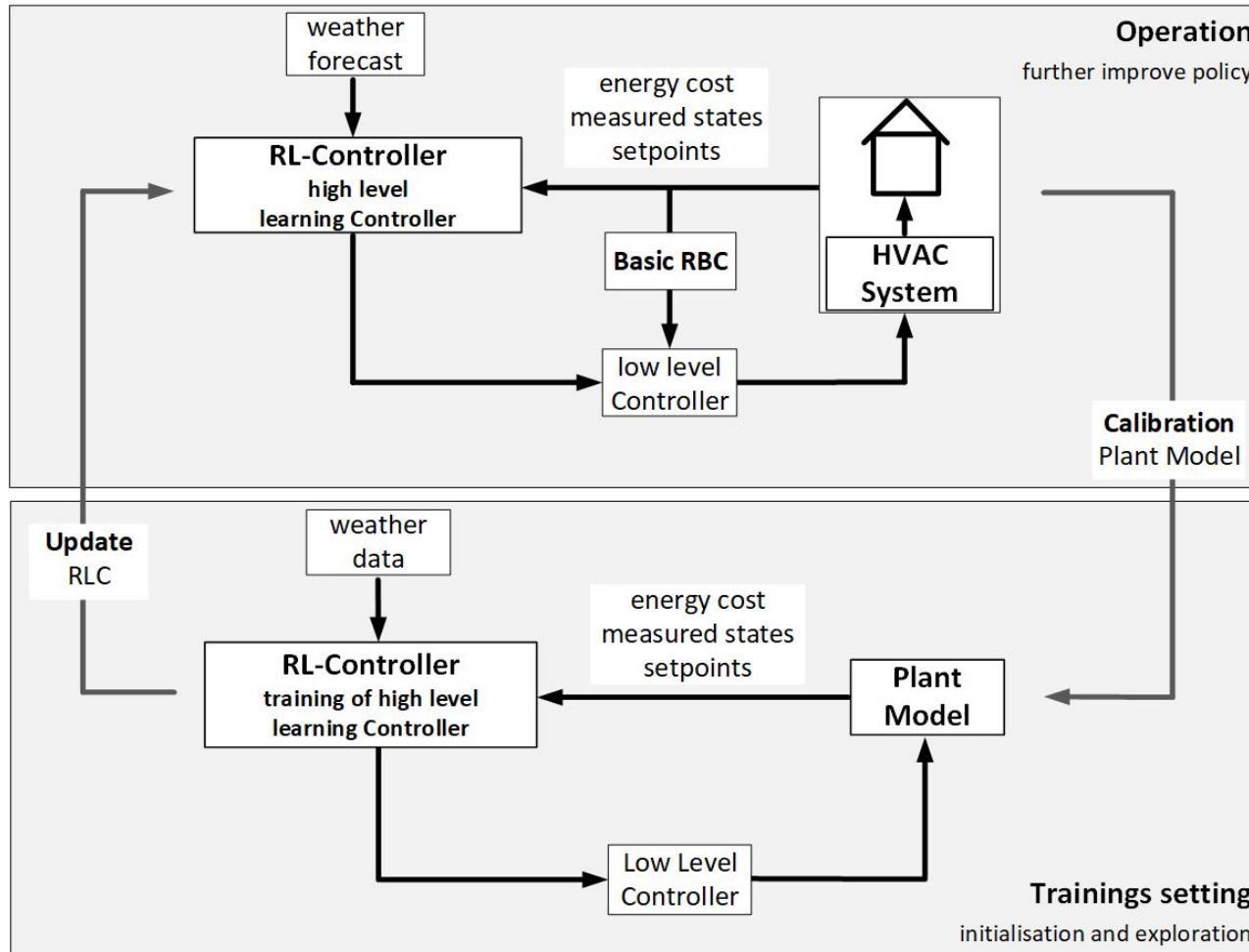
Reward (Energy and Discomfort Costs)



agent heats at
cheaper night rate
only

$-0.1 * (\text{energy rate})$

Simulated Reinforcement Learning



Schematic: Peter Bolt, ACSS IMAP ZHAW

RBC with Parameters Optimized by RL

- Some simple fixed rules
- Setpoints for water storage heating 1 & 2 are parametrized, k_1 and k_2 learned

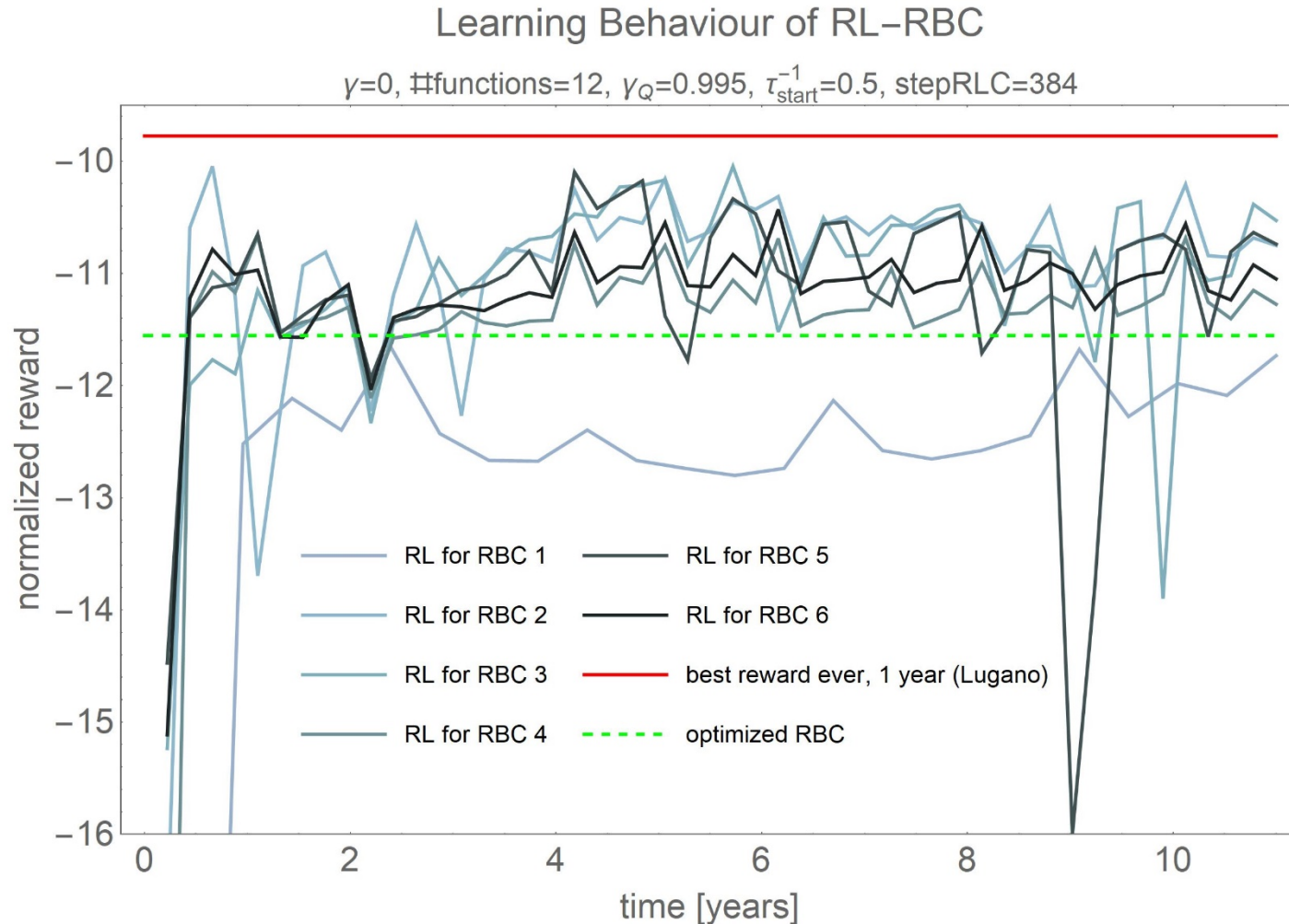
$$T_{sp,heating\ st1} = T_{sp,room} + \left(T_{sp,room} - \frac{T_{amb} + T_{forecast}}{2} \right) \cdot k_1$$

$$T_{sp,heating\ st2} = T_{sp,room} + \left(T_{sp,room} - \frac{T_{amb} + T_{forecast}}{2} \right) \cdot k_2$$

- Reward (energy consumption & set point violation) is normalized and thus almost independent of outside temperature:

$$R_{norm} = \frac{R}{\Delta T + 0.005\Delta T^2} \quad \text{with } \Delta T = T_{sp} - T_{amb}$$

RL with Parametrized RBC



→ acceptable performance in less than 1 year!

Summary

- RLC for a heating system
 - converges to nearly optimal trajectories
 - needs order of 100 simulated years for convergence (→ simulated RL)
- Improvements to RL
 - least-squares fit to get $q(s,a)$ in one step
 - n -step SARSA
 - truncated reward sum without discount
 - parallel, mutually supporting agents
- Alternative approach with parametrized RBC
 - much shorter learning time due to less coefficients
 - good performance