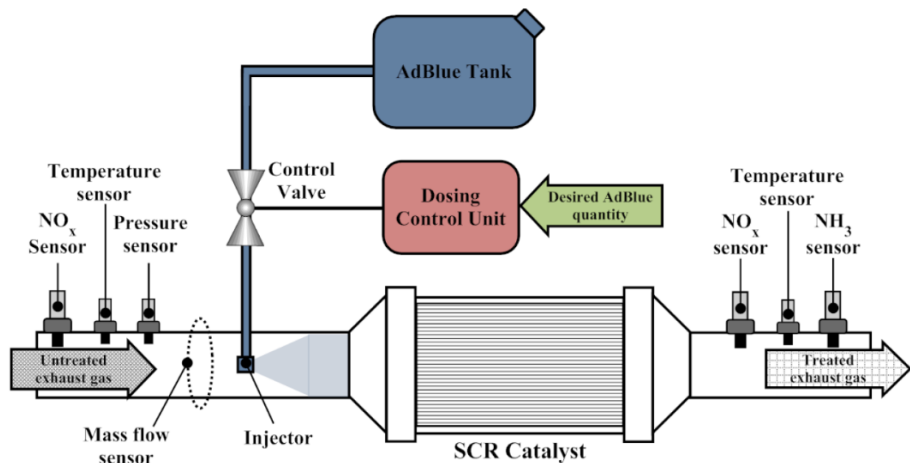# Systematic Catalytic Reduction by Reinforcement Learning

Chris Guiver, Baruch Gutierrez, Allen Hart, James Hook, Uwe Martin, Jordan Taylor

June 14, 2019
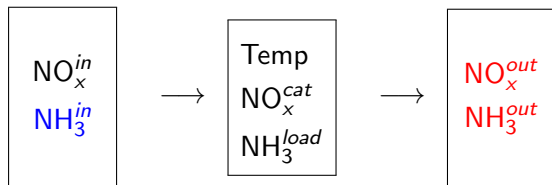
# Original Model

$$\frac{d}{dt}c_{NO,k} = \frac{n}{V_c \cdot \varepsilon_g} \cdot \frac{\overset{*}{m}_{EG} \cdot R}{p_{EG} \cdot M_{EG}} \left( T_{EG,k-1} \cdot c_{NO,k-1} - T_{c,k} \cdot c_{NO,k} \right) + $$
$$+ a_R \left( -4 \cdot r_{std,k} - 2 \cdot r_{fst,k} - r_{NO,g,k} \right)$$

$$\frac{d}{dt}c_{NO_2,k} = \frac{n}{V_c \cdot \varepsilon_g} \cdot \frac{\overset{*}{m}_{EG} \cdot R}{p_{EG} \cdot M_{EG}} \left( T_{EG,k-1} \cdot c_{NO_2,k-1} - T_{c,k} \cdot c_{NO_2,k} \right) + $$
$$+ a_R \left( -2 \cdot r_{fst,k} - 6 \cdot r_{slw,k} + r_{NO,g,k} \right)$$

$$\frac{d}{dt}c_{NH_3,k} = \frac{n}{V_c \cdot \varepsilon_g} \cdot \frac{\overset{*}{m}_{EG} \cdot R}{p_{EG} \cdot M_{EG}} \left( T_{EG,k-1} \cdot c_{NH_3,k-1} - T_{c,k} \cdot c_{NH_3,k} \right) + $$
$$+ a_R \left( -r_{ad,k} + r_{de,k} - 4 \cdot r_{ox,g,k} \right)$$

$$\frac{d}{dt}c_{O_2,k} = \frac{n}{V_c \cdot \varepsilon_g} \cdot \frac{\overset{*}{m}_{EG} \cdot R}{p_{EG} \cdot M_{EG}} \left( T_{EG,k-1} \cdot c_{O_2,k-1} - T_{c,k} \cdot c_{O_2,k} \right) + $$
$$+ a_R \left( -0.5 \cdot r_{NO,g,k} \right)$$

$$\frac{d}{dt}\theta_{NH_3,k} = \frac{1}{\Theta_{NH_3}} \left( r_{ad,k} - r_{de,k} - 4 \cdot r_{std,k} - 4 \cdot r_{fst,k} - 8 \cdot r_{slw,k} - 4 \cdot r_{ox,k} \right)$$

$$\frac{d}{dt}T_{c,k} = \frac{n}{m_c \cdot c_{p,c}} \left( \overset{*}{m}_{EG} \cdot c_{p,EG} \cdot \left( T_{EG,k-1} - T_{c,k} \right) + \alpha_c \cdot a_c \cdot \left( T_{Amb} - T_{c,k} \right) \right)$$

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathsf{T} = -a(\mathsf{T} - 200) + b\mathsf{NO}_x^{in}$$

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathsf{NO}_x^{cat} = \mathsf{NO}_x^{in} - \mathsf{NO}_x^{out} - \alpha R(\mathsf{T}, \mathsf{NO}_x^{cat}, \mathsf{NH}_3^{load})$$

$$\frac{\mathrm{d}}{\mathrm{d}t}\mathsf{NH}_3^{load} = \mathsf{NH}_3^{in} - \mathsf{NH}_3^{out} - R(\mathsf{T}, \mathsf{NO}_x^{cat}, \mathsf{NH}_3^{load})$$

$$R(\mathsf{T}, \mathsf{NO}_x^{cat}, \mathsf{NH}_3^{load}) = f(\mathsf{T})\mathsf{NO}_x^{cat}\mathsf{NH}_3^{load}$$

**Challenge:** Correctly control the amount of $NH_3^{in}$ injected into the system at every time step $t$.

**Challenge:** Correctly control the amount of $NH_3^{in}$ injected into the system at every time step $t$.

**Idea:** Use Reinforcement-Learning to train a controller capable of correctly deciding the optimal amount.

Reinforcement Learning involves an agent, moving through a state space, $\mathcal{S}$, by selecting an action from an action space, $\mathcal{A}$, at each state.

Given we are at some state $s_t \in \mathcal{S}$, taking an action $a_t \in \mathcal{A}$ provides the agent with the reward $r = r(s, a)$ and new state $s_{t+1}$.

The aim of the agent is to maximise the total (future) reward.

**SAMBa**

Idea:

$Q(s,a) \approx \mathbb{E} \left[ \begin{array}{l} \text{future discounted sum of rewards if we start at state} \\ \text{and then follow current policy for the rest of time} \end{array} \right]$
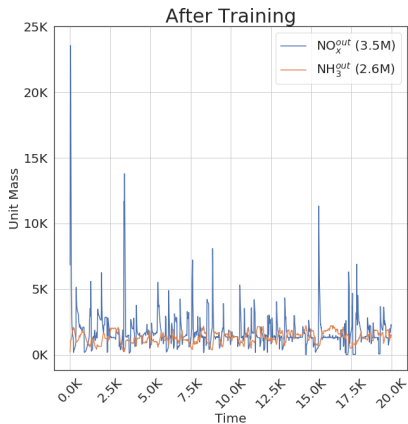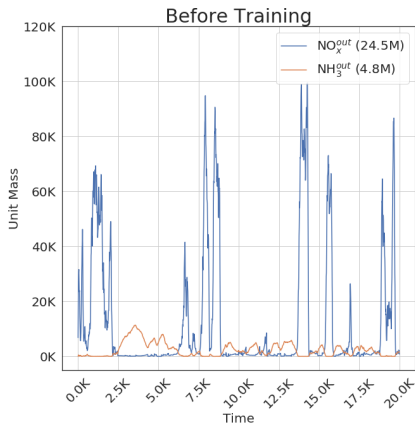
Ideally:

$$Q(s_t, a_t) = r_t + \gamma \max_a \mathbb{E}[Q(s_{t+1}, a)]$$

The objective of the training is to update $Q$ iteratively to take into account future values of $Q$, i.e. to correctly reflect the value of rewards available after multiple actions.

$$Q(s_t, a_t) \quad \longleftarrow \quad Q(s_t, a_t) - \alpha \left( Q(s_t, a_t) - r_t - \gamma \cdot \max_a \mathbb{E}[Q(s_{t+1}, a)] \right)$$
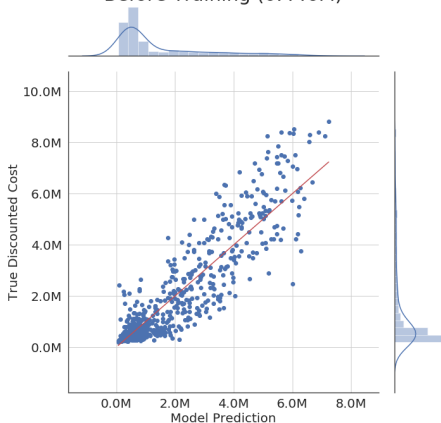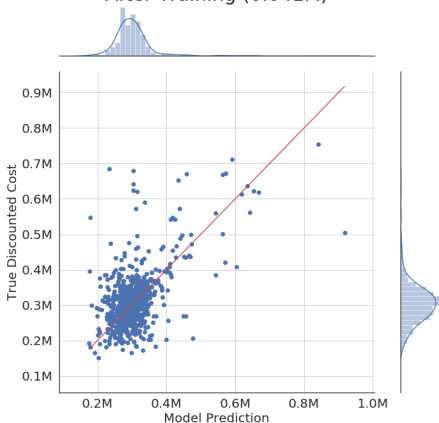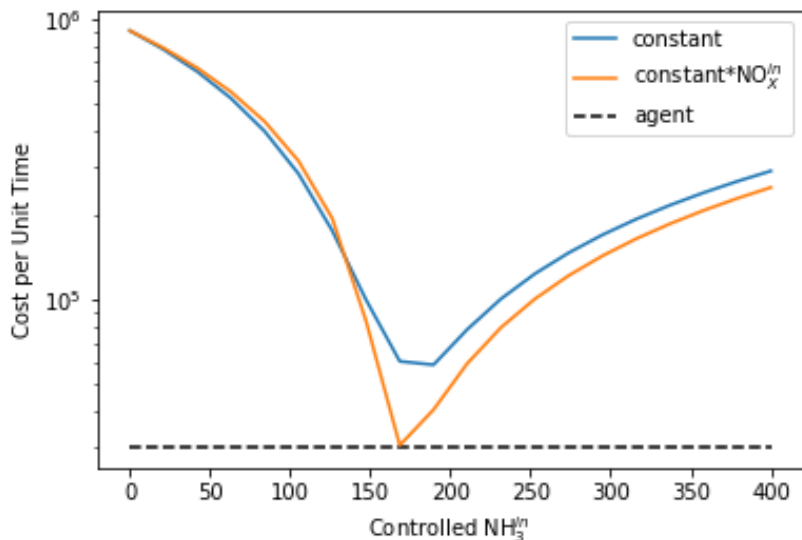
Before Training (0.446M)

After Training (0.041M)

Pros

- Our agent can learn a policy governed by a system of differential equations without seeing them.
- Ability to have "online" learning to cater policy to the user.
- Cheap evaluation to determine appropriate control.

Cons

- A parameter space to search i.e. discount factor $\gamma$, learning rate $\alpha$, and exploration rate $\epsilon$.
- Long training time.

Future

- Tune the toy model to be more realistic.
- Use the original set of differential equations.
- Allow noisy measurements.