

Automatic Generation Control for a Two Area Power System Using Deep Reinforcement Learning

Shane Reynolds

Charles Darwin University

Background

A generator's angular acceleration is governed by:

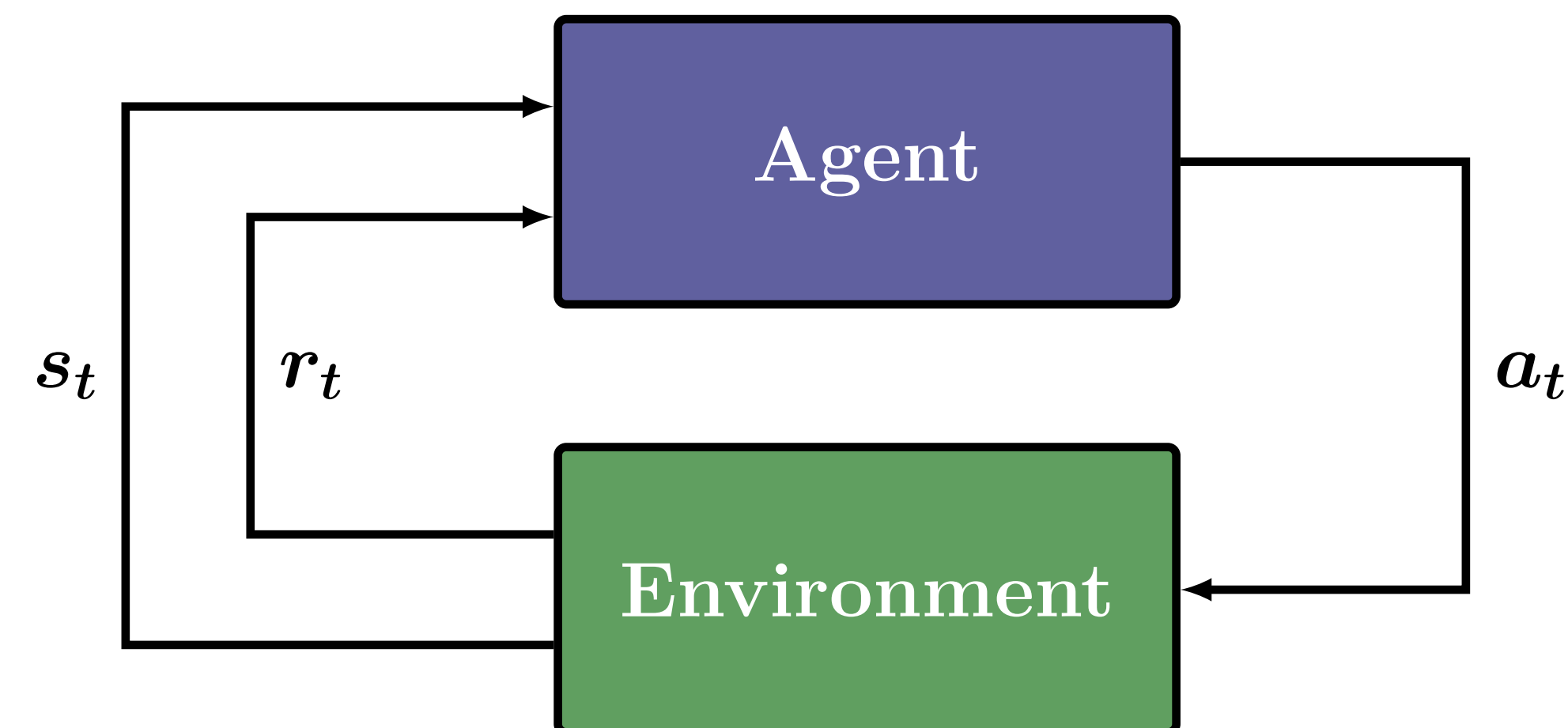
$$\Delta T = T_{mech} - T_{elec} = I\alpha$$

- If $\Delta T > 0$, then $\alpha \uparrow$ and f (Hz) \uparrow
- If $\Delta T < 0$, then $\alpha \downarrow$ and f (Hz) \downarrow

The Australian power network operates at 50 Hz.

Reinforcement Learning

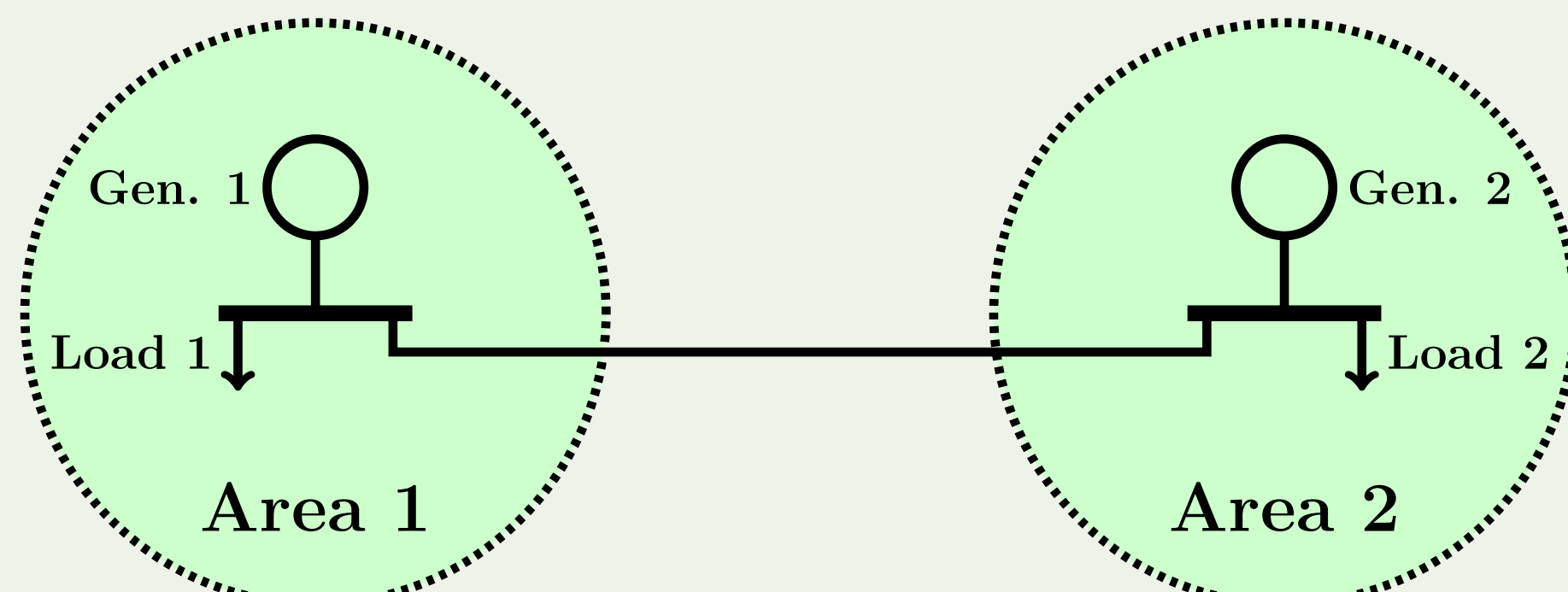
Reinforcement learning is a branch of machine learning concerned with an agent's sequential decision making to maximise cumulative expected reward.



The agent exists in some environment and at each time step observes state $s_t \in \mathcal{S}$; and takes an action $a_t \in \mathcal{A}$. Following this, the agent then receives a reward $r_t \in \mathcal{R} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [\mathcal{R}_{min}, \mathcal{R}_{max}]$.

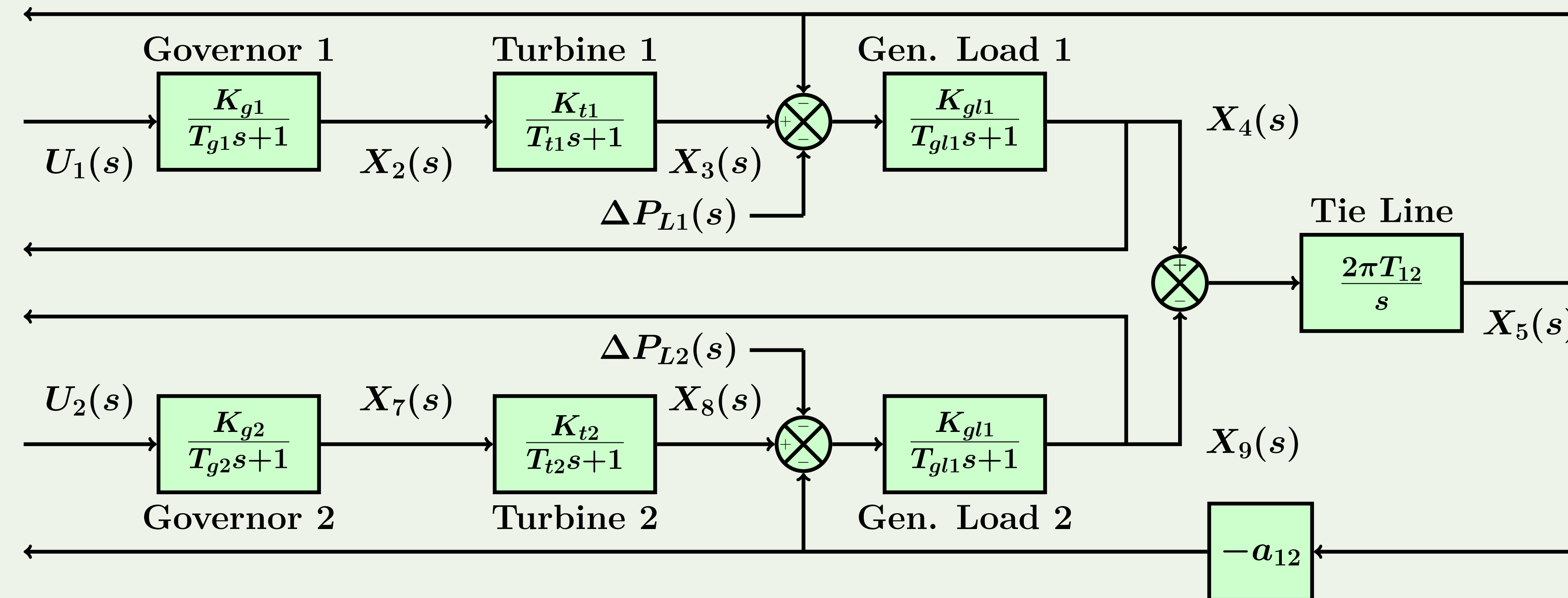
The Environment

Two power areas connected via a transmission line. Each power area consists of: a governor controlled generator; and stochastic load demand.



The control objective is to maintain inter-area power transfer, whilst regulating the frequency of each area.

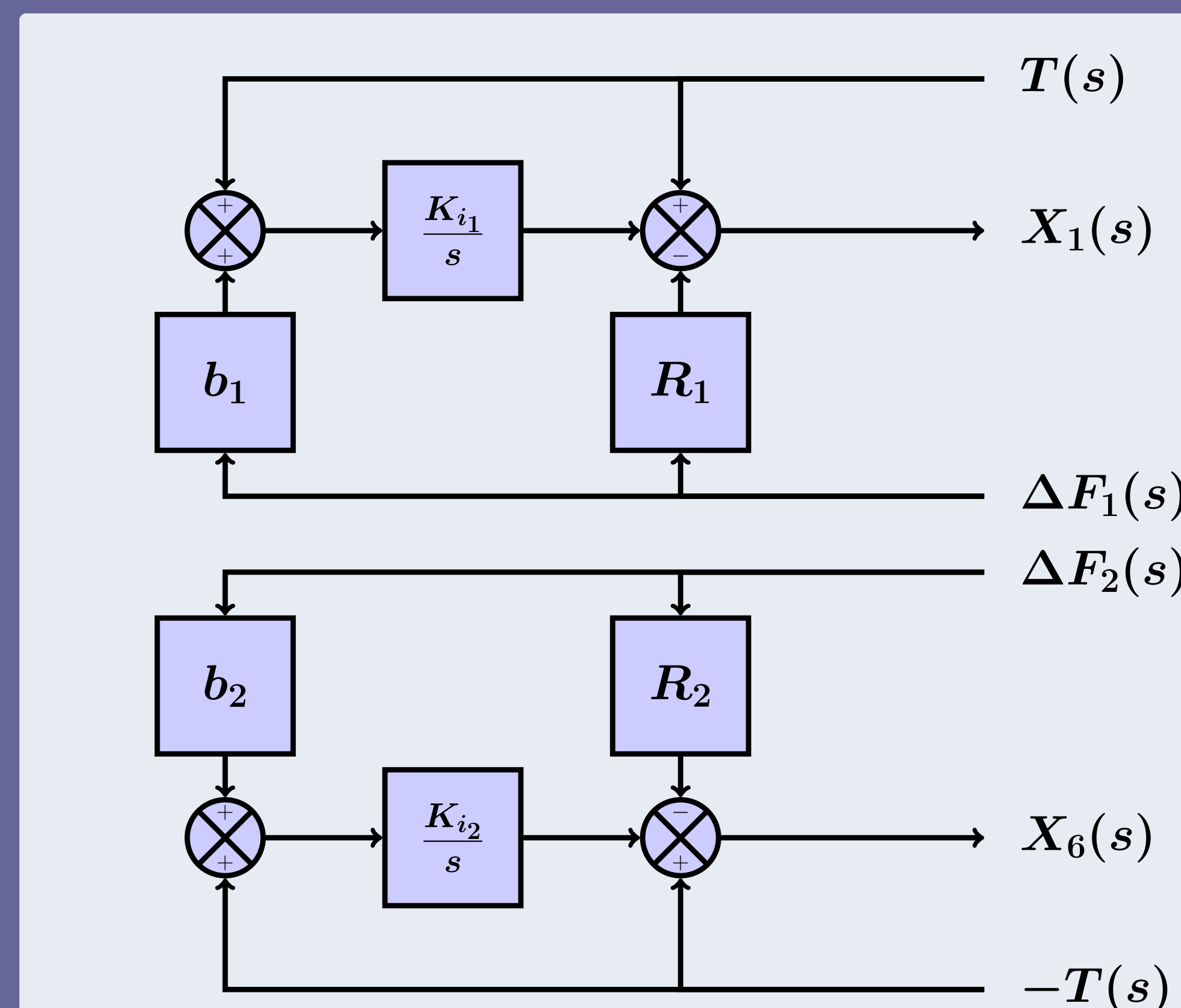
The Environment in the Frequency & Temporal Domains



Two Area System ODE

$$\begin{aligned} \dot{x}_2(t) &= \frac{1}{T_{sg1}} (K_{sg1} u_1(t) - x_2(t)) & \dot{x}_7(t) &= \frac{1}{T_{sg2}} (K_{sg2} u_2(t) - x_7(t)) \\ \dot{x}_3(t) &= \frac{1}{T_{t1}} (K_{t1} x_2(t) - x_3(t)) & \dot{x}_8(t) &= \frac{1}{T_{t2}} (K_{t2} x_7(t) - x_8(t)) \\ \dot{x}_4(t) &= \frac{1}{T_{gl1}} (K_{gl1} (x_3(t) - x_5(t) - \Delta p_{L1}(t)) - x_4(t)) & \dot{x}_9(t) &= \frac{1}{T_{gl2}} (K_{gl2} (x_8(t) - x_5(t) - \Delta p_{L2}(t)) - x_9(t)) \end{aligned}$$

Classical PI Controller



ODE System

$$\begin{aligned} \dot{x}_1(t) &= b_1 \Delta f_1(t) + x_5(t) \\ \dot{x}_6(t) &= b_2 \Delta f_2(t) - x_5(t) \end{aligned}$$

Experimental Setup

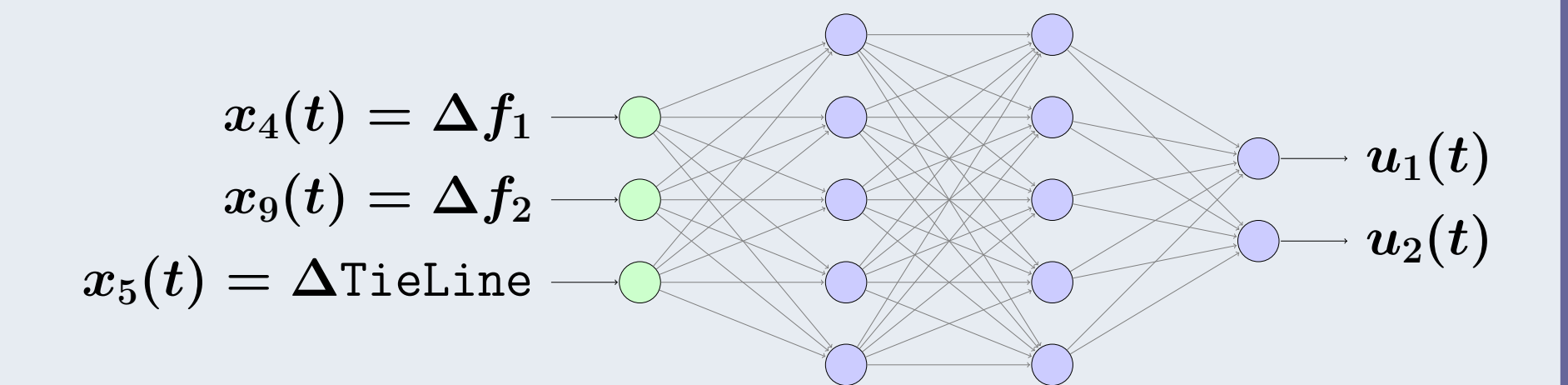
content...

Classical PI Controller Results

content...

DDPG Controller

Neural Network Architecture



Algorithm 1: DDPG

```

1 Rand. init. critic  $Q(s, a | \theta^Q)$  and actor  $\mu(s | \theta^\mu)$  with  $\theta^Q$  and  $\theta^\mu$ 
2 Init. target networks  $Q'$  and  $\mu'$  with  $\theta^{Q'} \leftarrow \theta^Q$ ,  $\theta^{\mu'} \leftarrow \theta^\mu$ 
3 Initialise replay buffer  $R$ 
4 for episode  $\leftarrow 1$  to  $M$  do
5   Initialise a random process  $\mathcal{N}$  for action exploration
6   Receive initial observation state  $s_1$ 
7   for  $t \leftarrow 0$  to  $T$  do
8     Select action  $a_t = \mu(s_t | \theta^\mu) + \mathcal{N}_t$  with noise exploration noise
9     Execute  $a_t$  and observe reward  $r_t$  and new state  $s_{t+1}$ 
10    Store transition  $(s_t, a_t, r_t, s_{t+1})$  in  $R$ 
11    Sample random  $N$  transitions  $(s_i, a_i, r_i, s_{i+1})$  from  $R$ 
12    Set  $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1} | \theta^{Q'}))$ 
13    Update critic by min. loss:  $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i | \theta^Q))^2$ 
14    Update the actor policy using the sampled policy gradient:
        
$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a | \theta^Q) |_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s | \theta^\mu) |_{s=s_i}$$

15    Update the target networks:
        
$$\theta^{Q'} = \tau \theta^Q + (1 - \tau) \theta^Q$$

        
$$\theta^{\mu'} = \tau \theta^\mu + (1 - \tau) \theta^\mu$$

16 end
17 end

```

DDPG Controller Results