# Switching-time bioprocess control with pulse-width-modulated optogenetics

**Sebastián Espinel-Ríos** *

* *School of Chemical and Bioprocess Engineering, University College Dublin, Ireland (Starting January 1, 2026)*

**Abstract:** Biotechnology can benefit from dynamic control to improve production efficiency. In this context, optogenetics enables modulation of gene expression using light as an external input, allowing fine-tuning of protein levels to unlock dynamic metabolic control and regulation of cell growth. Optogenetic systems can be actuated by light intensity. However, relying solely on intensity-driven control (i.e., signal amplitude) may fail to properly tune optogenetic bioprocesses when the dose-response relationship (i.e., light intensity versus gene-expression strength) is steep. In these cases, tunability is effectively constrained to either fully active or fully repressed gene expression, with little intermediate regulation. Pulse-width modulation, a concept widely used in electronics, can alleviate this issue by alternating between fully ON and OFF light intensity within forcing periods, thereby smoothing the average response and enhancing process controllability. Naturally, optimizing pulse-width-modulated optogenetics entails a switching-time optimal control problem with a binary input over many forcing periods. While this can be formulated as a mixed-integer program on a refined time grid, the number of decision variables can grow rapidly with increasing time-grid resolution and number of forcing periods, compromising tractability. Here, we propose an alternative solution based on reinforcement learning. We parametrize control actions via the duty cycle, a continuous variable that encodes the ON-to-OFF switching time within each forcing period, thereby respecting the intrinsic binary nature of the light intensity.

*Keywords:* Bioprocess control, switching-time control, pulse-width modulation, duty cycles, reinforcement learning, uncertainty.

## 1. INTRODUCTION

Biotechnology leverages the enzymatic pathways of microorganisms to synthesize chemicals, materials, and fuels, among other valuable products. Microbial bioproduction therefore plays a crucial role in the development of a future circular economy (Ewing et al., 2022; Konzock and Nielsen, 2024). One emerging technology to optimize bioproduction is optogenetics, which can be used to fine-tune gene expression in real-time using light as an external control input (Milias-Argeitis et al., 2016; Hoffman et al., 2022; Benisch et al., 2024). Either via gene expression activation or repression, this enables dynamic modulation of intracellular enzyme levels throughout bioprocesses, for example, for dynamic metabolic control. Optogenetics has also shown potential for regulating microbial growth through modulation of the expression of proteins involved in toxin-antitoxin systems, antibiotic-resistance modules, or auxotrophic amino acid synthesis. These capabilities promise enhanced control of population levels in synthetic microbial communities.

Optogenetic bioprocesses can be actuated via light intensity (i.e., a continuous signal amplitude) of specific wavelengths. Light-intensity-driven optogenetics has been explored through model-based predictive control and reinforcement learning (RL) (e.g., Milias-Argeitis et al. (2016); Espinel-Ríos et al. (2025a,b)). In an optimal control context, light intensity can be discretized over finite piecewise-constant control intervals that act as the degrees of freedom. However, intensity-driven optogenetic control can suffer from limited tunability and robustness (Davidson et al., 2013; Benzinger et al., 2022). In many optogenetic systems, dose-response curves (i.e., light intensity versus gene activation/repression) are very steep, driving the system close to fully active or fully repressed gene-expression states and thus hindering intermediate operating points.

Pulse-width modulation (PWM) is an alternative way to drive optogenetic bioprocesses (Davidson et al., 2013; Benzinger et al., 2022) that can address the issue of steep dose-response. It exploits two light-intensity levels (ON/OFF) over predefined forcing periods (a full interval of ON and OFF subintervals). The ON state typically corresponds to the maximum achievable light intensity, and the OFF state to zero intensity. This strategy improves practical tunability and robustness by *averaging* the response over the forcing period.

Optimizing PWM-driven bioprocesses is, however, non-trivial. The optimal solution of a PWM-driven optogenetic bioprocess entails optimizing the switching times for the many forcing periods throughout the process. One way to formulate this is to discretize time and optimize constant binary inputs within each forcing period while enforcing the constraint that once an input switches OFF, all remaining inputs in that period must remain OFF, thereby preserving the sequential ON-OFF pattern. This

leads to a *switching-time optimal control problem with a binary input over multiple forcing periods*. Solving this as a *mixed-integer programming problem*, at first glance a sound option, would require extremely fine control-interval discretization to accurately approximate the continuous switching points. Yet, such formulations quickly become computationally intractable as the time grid becomes finer and the number of forcing periods grows. Coarser discretization may facilitate tractability but restricts the solution space. Here, we use RL with duty cycles as continuous decision variables that encode binary light inputs. RL policies can be pretrained *in silico* (e.g., using a digital twin).

The remainder of this paper is structured as follows. Section 2 formulates the switching-time optimal control problem underlying PWM. Section 3 outlines our proposed strategy based on RL. Finally, Section 4 presents a case study involving optogenetic control of *Escherichia coli* growth via PWM. Control policy robustness is demonstrated via randomization on the dynamics and initial conditions.

## 2. SWITCHING-TIME OPTIMAL CONTROL

In PWM-driven optogenetics, the optogenetic input $\boldsymbol{u}(t)$ is modeled as a binary-valued vector of dimension $n_u$, where $n_u$ is the number of inputs (i.e., independent light intensity channels). For each input channel $i$, we define:

$$u_i(t) \in \{0, 1\}, \tag{1}$$

where $u_i(t) = 0$ (OFF) encodes $0\,\%$ and $u_i(t) = 1$ (ON) encodes $100\,\%$ of a given maximum light intensity $I_{i,\max} \in \mathbb{R}$. The corresponding *physical* light intensity for channel $i$ is thus:

$$I_i(u_i(t)) = u_i(t)\, I_{i,\max}. \tag{2}$$

Furthermore, the overall process involves $n_T$ forcing periods, so that $t_f = n_T T$ denotes the finite process time. We define the $k$-th forcing period $\mathcal{T}_k$ as the interval:

$$\mathcal{T}_k := [kT, (k+1)T], \quad \forall k \in \{0, 1, \dots, n_T - 1\}. \tag{3}$$

The trajectory of input $i$ over the $n_T$ forcing periods of the process follows a binary ON-OFF pattern within each forcing period:

$$u_{i,k}(t) = \begin{cases} 1, & t \in [kT, \tau_{i,k}), \\ 0, & t \in [\tau_{i,k}, (k+1)T), \end{cases} \tag{4}$$
$$\forall k \in \{0, 1, \dots, n_T - 1\},$$

where $kT \leq \tau_{i,k} \leq (k+1)T$ denotes the switching time at which the input transitions from ON to OFF within the period $\mathcal{T}_k$. In a PWM-driven optogenetic bioprocess, the decision variable within the period $\mathcal{T}_k$ is therefore the switching time $\tau_{i,k}$. Here we assume that each forcing period has fixed duration $T$, known *a priori*.

Let us now consider an optogenetic bioprocess with state $\boldsymbol{x}(t) \in \mathbb{R}^{n_x}$ and measured output $\boldsymbol{y}(t) = \boldsymbol{h}(\boldsymbol{x}(t))$, where $\boldsymbol{h} : \mathbb{R}^{n_x} \mapsto \mathbb{R}^{n_y}$ captures the underlying measurement function. The continuous-time dynamics read:

$$\frac{\mathrm{d}\boldsymbol{x}(t)}{\mathrm{d}t} = f\left(\boldsymbol{x}(t), \boldsymbol{I}(t)\right), \quad \boldsymbol{x}(t_0) = \boldsymbol{x}_0, \tag{5}$$

where $f : \mathbb{R}^{n_x} \times \mathbb{R}^{n_u} \mapsto \mathbb{R}^{n_x}$ captures the underlying process dynamics, $\boldsymbol{I}(t) \in \mathbb{R}^{n_u}$ is the vector of physical

light intensities, $t_0$ is the initial time, and $\boldsymbol{x}_0$ the initial state.

With these definitions in place, the generalized *switching-time optimal control problem with a binary input* is:

$$\min_{\{\tau_{i,k}\}_{i=1,\ k=0}^{n_u,\ n_T-1}} \mathbb{E}[J(\cdot)], \tag{6a}$$

$$\text{s.t.} \qquad \frac{\mathrm{d}\boldsymbol{x}(t)}{\mathrm{d}t} = f\left(\boldsymbol{x}(t), \boldsymbol{I}(t)\right), \tag{6b}$$

$$\boldsymbol{x}(t_0) = \boldsymbol{x}_0, \tag{6c}$$

$$\boldsymbol{y}(t) = \boldsymbol{h}(\boldsymbol{x}(t)), \tag{6d}$$

$$u_{i,k}(t) = \begin{cases} 1, & t \in [kT, \tau_{i,k}), \\ 0, & t \in [\tau_{i,k}, (k+1)T), \end{cases} \tag{6e}$$

$$I_{i,k}(u_{i,k}(t)) = u_{i,k}(t)\, I_{i,\max}, \tag{6f}$$

$$kT \leq \tau_{i,k} \leq (k+1)T, \tag{6g}$$

$$\forall k \in \{0, 1, \dots, n_T - 1\},$$
$$\forall i \in \{1, \dots, n_u\}.$$

Here, $\mathbb{E}[J(\cdot)]$ denotes the expectation of an optimization objective function $J(\cdot)$, and the decision variables are the switching times $\tau_{i,k}$ for all inputs across all forcing periods.

Without loss of generality, in this paper we consider reference tracking control:

$$J = \int_{t_0}^{t_f} \|\boldsymbol{y}(t) - \boldsymbol{r}(t)\|_{\mathbf{Q_s}}^2\, dt + \|\boldsymbol{y}(t_f) - \boldsymbol{r}(t_f)\|_{\mathbf{Q_t}}^2, \tag{7}$$

where the reference trajectory for the controlled output variable $\boldsymbol{y}(t)$ is denoted as $\boldsymbol{r}(t)$. $\mathbf{Q_s}$ and $\mathbf{Q_t}$ are weight matrices of appropriate dimension, weighting the contribution of the *stage* tracking cost and the *terminal* tracking cost, respectively. In this notation, $\|\boldsymbol{a}\|_{\mathbf{A}}^2 := \boldsymbol{a}^\top \mathbf{A} \boldsymbol{a}$ denotes the squared norm of a vector $\boldsymbol{a}$ weighted by the matrix $\mathbf{A}$.

## 3. SOLUTION STRATEGY

As mentioned in the Introduction, Problem (6) can in principle be viewed through the lens of *mixed-integer programming*, involving binary inputs and continuous-time dynamics discretized on a time grid. However, such a formulation can quickly become intractable due to the *curse of dimensionality* in the number of decision variables.

Here, we instead follow an RL-based solution approach. The key idea is *not* to optimize the piecewise binary inputs over a fine time grid explicitly (e.g., as in mixed-integer programming), but rather to work with duty cycles as continuous decision variables that *implicitly* encode the underlying binary inputs, leading to a more tractable *continuous* optimization problem.

In a forcing period $\mathcal{T}_k$, the duty cycle of input $i$ (cf. Eq. (4)), $D_{i,k} \in [0, 1]$, is a continuous variable defined by:

$$D_{i,k} = \frac{1}{T} \int_{kT}^{(k+1)T} u_{i,k}(t)\, dt = \frac{1}{T} \int_{kT}^{\tau_{i,k}} 1\, dt = \frac{\tau_{i,k} - kT}{T},$$
$$\forall k \in \{0, 1, \dots, n_T - 1\}, \forall i \in \{1, \dots, n_u\}. \tag{8}$$

As a consequence, $\tau_{i,k}$ is uniquely determined by $D_{i,k}$ via:

$$\tau_{i,k}(D_{i,k}) = \left(k + D_{i,k}\right)T. \tag{9}$$

Let $\boldsymbol{D}_k \in [0,1]^{n_u}$ collect the duty cycles of all inputs in forcing period $k$, $\boldsymbol{\tau}_k \in [kT,(k+1)T]^{n_u}$ the corresponding switching times, $\boldsymbol{u}_k(t) := [u_{1,k}(t),\ldots,u_{n_u,k}(t)]^\top \in \{0,1\}^{n_u}$ the time-varying binary input vector over $\mathcal{T}_k$, and $\boldsymbol{I}_k(t) := [I_{1,k}(t),\ldots,I_{n_u,k}(t)]^\top \in \mathbb{R}^{n_u}$ the corresponding physical light intensity. We treat $\boldsymbol{D}_k$ as the decision variable over $\mathcal{T}_k$ since $\boldsymbol{D}_k \mapsto \boldsymbol{\tau}_k(\boldsymbol{D}_k) \mapsto \boldsymbol{u}_k(t) \mapsto \boldsymbol{I}_k(t)$, which ultimately drives the system dynamics.

Furthermore, we describe the state transition (cf. Eq. (5)) as a Markov decision process:

$$\boldsymbol{x}_{k+1} \sim \mathrm{P}\big(\boldsymbol{x}_{k+1} \mid \boldsymbol{x}_k, \boldsymbol{D}_k\big),\ k \in \{0,1,\ldots,n_T-1\}, \quad (10)$$

where P denotes the conditional probability.

Let $\pi(\boldsymbol{D}_k \mid \boldsymbol{x}_k, \boldsymbol{\theta})$ be the control policy, parametrized by $\boldsymbol{\theta} \in \mathbb{R}^{n_\theta}$. We transform the optimization problem (6) to:

$$\max_{\boldsymbol{\theta}}\ \mathbb{E}_{\boldsymbol{\tau} \sim \mathrm{P}(\boldsymbol{\tau}|\boldsymbol{\theta})}\left[J(\boldsymbol{\tau})\right], \quad (11)$$

where $\boldsymbol{\tau}$ is a trajectory of states, inputs, rewards, and transitions:

$$\boldsymbol{\tau} = \{(\boldsymbol{x}_k, \boldsymbol{D}_k, r_{k+1}, \boldsymbol{x}_{k+1})\}_{k=0}^{n_T-1}, \quad (12)$$

and the reward at time $k$ is defined as:

$$r_k := -\|\boldsymbol{y}_k - \boldsymbol{r}_k\|_{\mathbf{Q_s}}^2, \quad k \in \{1,\ldots,n_T-1\}, \quad (13a)$$

$$r_{n_T} := -\|\boldsymbol{y}_{n_T} - \boldsymbol{r}_{n_T}\|_{\mathbf{Q_t}}^2. \quad (13b)$$

upon discretizing the objective function in (7) over the forcing periods. Note that we collect the first reward after the first duty-cycle input has been applied, consistent with the definition in Eq. (12).

The return in this RL formulation is then defined as:

$$J := \sum_{k=1}^{n_T} r_k, \quad (14)$$

which corresponds to an *undiscounted episodic return*. It can be seen as the discrete-time counterpart of the objective in Eq. (7), using negative costs so that maximizing $J$ in (11) is equivalent to minimizing Eq. (7).

To solve (11), we apply policy gradients in a gradient-ascent fashion over $n_m$ epochs, with learning rate $\alpha$:

$$\boldsymbol{\theta}_{m+1} = \boldsymbol{\theta}_m + \alpha \nabla_{\boldsymbol{\theta}} \mathbb{E}_{\boldsymbol{\tau}}\left[J(\boldsymbol{\tau})\right], \quad m = 0,\ldots,n_m-2. \quad (15)$$

The joint probability of the trajectory $\boldsymbol{\tau}$ follows:

$$\mathrm{P}(\boldsymbol{\tau} \mid \boldsymbol{\theta}) = \mathrm{P}(\boldsymbol{x}_0) \prod_{k=0}^{n_T-1} \left[\pi(\boldsymbol{D}_k \mid \boldsymbol{x}_k, \boldsymbol{\theta})\,\mathrm{P}(\boldsymbol{x}_{k+1} \mid \boldsymbol{x}_k, \boldsymbol{D}_k)\right]. \quad (16)$$

Applying the Policy Gradient Theorem (Sutton et al., 1999) and approximating the expectation via $n_{\mathrm{MC}}$ Monte Carlo simulations leads to:

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{\boldsymbol{\tau}}\left[J(\boldsymbol{\tau})\right] = \mathbb{E}_{\boldsymbol{\tau}}\left[J(\boldsymbol{\tau}) \nabla_{\boldsymbol{\theta}}\left[\sum_{k=0}^{n_T-1} \log \pi(\boldsymbol{D}_k \mid \boldsymbol{x}_k, \boldsymbol{\theta})\right]\right],$$

$$\approx \frac{1}{n_{\mathrm{MC}}} \sum_{j=1}^{n_{\mathrm{MC}}} \left[\left(\frac{J(\boldsymbol{\tau}^{(j)}) - \bar{J}_m}{\sigma_{J_m} + \epsilon}\right)\right.$$

$$\left.\cdot \nabla_{\boldsymbol{\theta}}\left[\sum_{k=0}^{n_T-1} \log\left(\pi(\boldsymbol{D}_k^{(j)} \mid \boldsymbol{x}_k^{(j)}, \boldsymbol{\theta})\right)\right]\right], \quad (17)$$

with a baseline normalized by the mean $\bar{J}_m$ and the standard deviation $\sigma_{J_m}$ of the return in epoch $m$; $\epsilon$ is a small positive scalar.

# 4. PWM-DRIVEN CONTROL OF OPTOGENETIC CELL GROWTH

We consider an *E. coli* strain with an engineered lysine auxotrophy obtained by deletion of *lysA* (encoding diaminopimelate decarboxylase), which is essential for lysine biosynthesis.

## 4.1 Mathematical model

The corresponding macroscopic bioprocess dynamics for optogenetic growth in a well-mixed chemostat with dilution rate $d_l$ follow (Espinel-Ríos et al., 2025a):

$$\frac{\mathrm{d}b}{\mathrm{d}t} = \big(\mu(g,p) - d_l\big)\,b,\ b(0) = b_0, \quad (18a)$$

$$\frac{\mathrm{d}g}{\mathrm{d}t} = -q_g(g,p)\,b + (g_{\mathrm{in}} - g)\,d_l,\ g(0) = g_0, \quad (18b)$$

$$\frac{\mathrm{d}p}{\mathrm{d}t} = q_p(I) - \big(d_p + \mu(g,p)\big)\,p,\ p(0) = p_0, \quad (18c)$$

with kinetic rates:

$$\mu(g,p) = \mu_{\max}\left(\frac{g}{g+k_g}\right)\left(\frac{f_c p}{f_c p + k_p}\right), \quad (19a)$$

$$q_g(g,p) = Y_{g/b}\,\mu(g,p), \quad (19b)$$

$$q_p(I) = q_{p,\max}\left(\frac{I^n}{I^n + k_I^n}\right), \quad (19c)$$

where $b \in \mathbb{R}$ is the biomass concentration, $g \in \mathbb{R}$ the extracellular glucose concentration, and $p \in \mathbb{R}$ the intracellular lysine concentration. $I \in \mathbb{R}$ is the light intensity. The intensity-dependent lysine synthesis rate $q_p(I) : \mathbb{R} \mapsto \mathbb{R}$ is a Hill function that lumps the light-controlled expression of *lysA* together with the resulting lysine production. The parameter $d_p$ accounts for lysine conversion and usage to support growth and viability.

We assume a ccaS/ccaR two-component optogenetic system (Davidson et al., 2013), which is inducible by green light. We assume the following parameters: $n = 0.2191$ and $k_I = 5.5086 \times 10^{-7}\,\mathrm{W/m^2}$, adapted and fitted (without leakage) from the dose-response data in Davidson et al. (2013), and $q_{p,\max} = 0.3366\,\mathrm{mmol/(g \cdot h)}$ (from Espinel-Ríos et al. (2025a)). Light intensity is measured in $\mathrm{W/m^2}$. The remaining parameters are (Espinel-Ríos et al., 2025a): $\mu_{\max} = 0.982\,\mathrm{h^{-1}}$, $f_c = 1100\,\mathrm{g/L}$, $Y_{g/b} = 10.18\,\mathrm{mmol/g}$, $k_g = 2.964 \times 10^{-4}\,\mathrm{mmol/L}$, $k_p = 1.7\,\mathrm{mmol/L}$, $d_l = 0.15\,\mathrm{h^{-1}}$, $g_{\mathrm{in}} = 200\,\mathrm{mmol/L}$, $d_p = 20.8\,\mathrm{h^{-1}}$.

## 4.2 Smoothing dose-response curves via duty cycles

Before diving into the implementation of the RL strategy, we demonstrate the smoothing effect of PWM on steep dose-response curves via duty-cycle averaging. Let us consider a binary ON-OFF input sequence over the forcing period $\mathcal{T}_k$ and one light-intensity channel (green light) ($n_u = 1$), with $I_{1,\max} = 30\,\mathrm{W/m^2}$ and duty cycle $D_{1,k} \in [0,1]$. Hereafter, we omit the subscript indicating input channel 1 for simplicity, as there is only one input channel in the process ($I_{1,\max} := I_{\max}$, $D_{1,k} := D_k$). During the ON subinterval, the intensity is $I(t) = I_{\max}$, and during the OFF subinterval it is $I(t) = 0$ (cf. Eq. (2)).

The period-averaged value of the Hill function, denoted as $\bar{q}_p$, over the interval $[kT,(k+1)T]$ is:
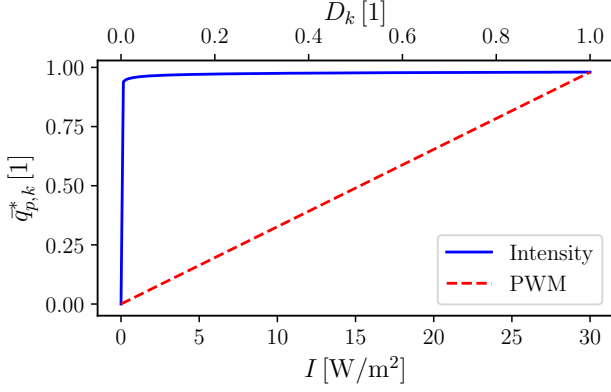
Fig. 1. Comparison of the *normalized* average Hill activation function $\bar{q}_{p,k}^*$ over the forcing period $\mathcal{T}_k$ under intensity-driven and PWM-driven actuation. The normalized average activation is defined as $\bar{q}_{p,k}^* := \bar{q}_{p,k}/q_{p,\max}$, yielding a range $[0, 1]$.

$$\bar{q}_{p,k} = \frac{1}{T}\int_{kT}^{(k+1)T} q_p(I(t))\,\mathrm{d}t$$
$$= \frac{1}{T}\left[\int_{kT}^{\tau_k=(k+D_k)T} q_p(I_{\max})\,\mathrm{d}t + \int_{\tau_k=(k+D_k)T}^{(k+1)T} q_p(0)\,\mathrm{d}t\right]$$
$$= D_k\, q_p(I_{\max}). \tag{20}$$

Thus, it becomes clear that the average Hill activation over the period is determined solely by the duty cycle $D_k$.

Note that the parameters of the Hill function introduced in the previous section yield a very steep activation curve. Fig. 1 shows that duty-cycle modulation results in a much smoother effective input-output relationship than direct intensity control, facilitating practical tunability.

### 4.3 RL-derived control policies

To evaluate the proposed RL-based solution to the switching-time optimal control problem, we consider a three-setpoint biomass reference trajectory involving transitions from $3\,\mathrm{g/L}$ to $5\,\mathrm{g/L}$ and finally to $7\,\mathrm{g/L}$ over $24\,\mathrm{h}$, with 24 forcing periods of $1\,\mathrm{h}$ each. Such a trajectory may arise in synthetic microbial consortia where biomass concentration is linked to the productivity of a specific metabolic submodule carried by a specific strain. The corresponding control problem is to find the 24 duty cycles (encoding the ON-to-OFF switching times) for the 24 forcing periods that best track this reference trajectory.

We now introduce our control scenarios:

- **Ideal system**: we assume no uncertainty in the process dynamics.
- **Uncertain system**: we assume different uncertainty levels in the process dynamics. Specifically, we assume that the initial conditions of the plant and the maximum gene expression rate $q_{p,\max}$ follow normal distributions with standard deviations of $2.5, 5,$ and $7.5\,\%$ relative to their mean (nominal) values.

We use the model in Section 4.1 as a *digital twin* for our optogenetic bioprocess, and thus as the virtual environment to train our RL policies. The nominal initial conditions

are $b_0 = 3\,\mathrm{g/L}$, $g_0 = 50\,\mathrm{mmol/L}$, and $p_0 = 1.0752 \times 10^{-4}\,\mathrm{mmol/g}$.
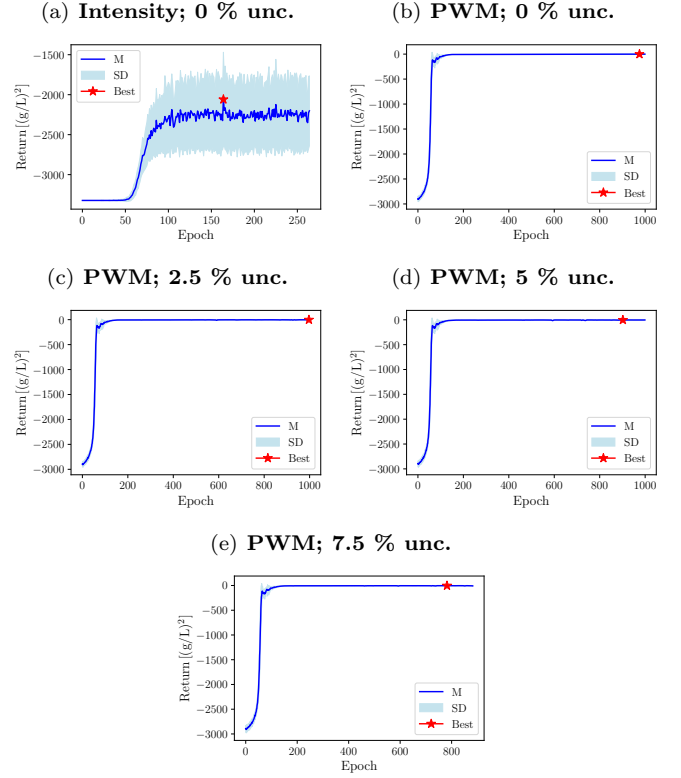


Fig. 2. Return over training epochs for the RL policies under intensity-driven and PWM-driven optogenetic control for selected uncertainty (unc.) levels. M: mean; SD: standard deviation; Best: selected policy.

We train RL policies as outlined in Section 3 in PyTorch (Paszke et al., 2019). We use a fully connected stochastic policy network (4 hidden layers, 20 neurons each, Leaky ReLU). The output linear layer returns the mean and standard deviation of the duty cycle, following a normal distribution. We use $n_m = 1000$, $n_{\mathrm{MC}} = 100$, and $\alpha = 0.001$. Early stopping is applied if no improvement in the return is observed for 100 consecutive epochs. To facilitate convergence, we augment the policy input with the current and previous states and duty cycles, as well as a process-time embedding $e_k \in [-1, 1]$. Accordingly, the policy is reformulated as $\pi\big(D_k \mid \boldsymbol{x}_k, \boldsymbol{\theta}\big) := \pi\big(D_k \mid \boldsymbol{z}_k, \boldsymbol{\theta}\big)$, where $\boldsymbol{z}_k := \big[\boldsymbol{x}_k^\top,\ D_{k-1},\ \boldsymbol{x}_{k-1}^\top,\ D_{k-2},\ e_k\big]^\top$.

Fig. 2 shows the evolution of the return over training epochs, together with the selected (best) policy. As a benchmark, we also trained a policy based directly on light intensity: here, the intensity was allowed to vary in the range $[0, I_{\max}]$ as a piecewise-constant input with the same interval length as the forcing periods in the PWM cases. Even in the absence of system uncertainty, the intensity-driven control scenario resulted in very low returns with large variability, as reflected by the large standard deviation of the return (cf. Fig. 2a). This can be explained by the extremely steep Hill-type activation $q_p(I)$ (cf. Fig. 1), which induces almost binary ON-OFF gene-expression transitions and thus yields weakly informative gradients for learning. This effectively stalls the learning process, even though *theoretically* speaking intensity actu-

**(a) Intensity; 0 % uncertainty**

**(b) PWM; 0 % uncertainty**

**(c) PWM; 2.5 % uncertainty**

**(d) PWM; 5 % uncertainty**
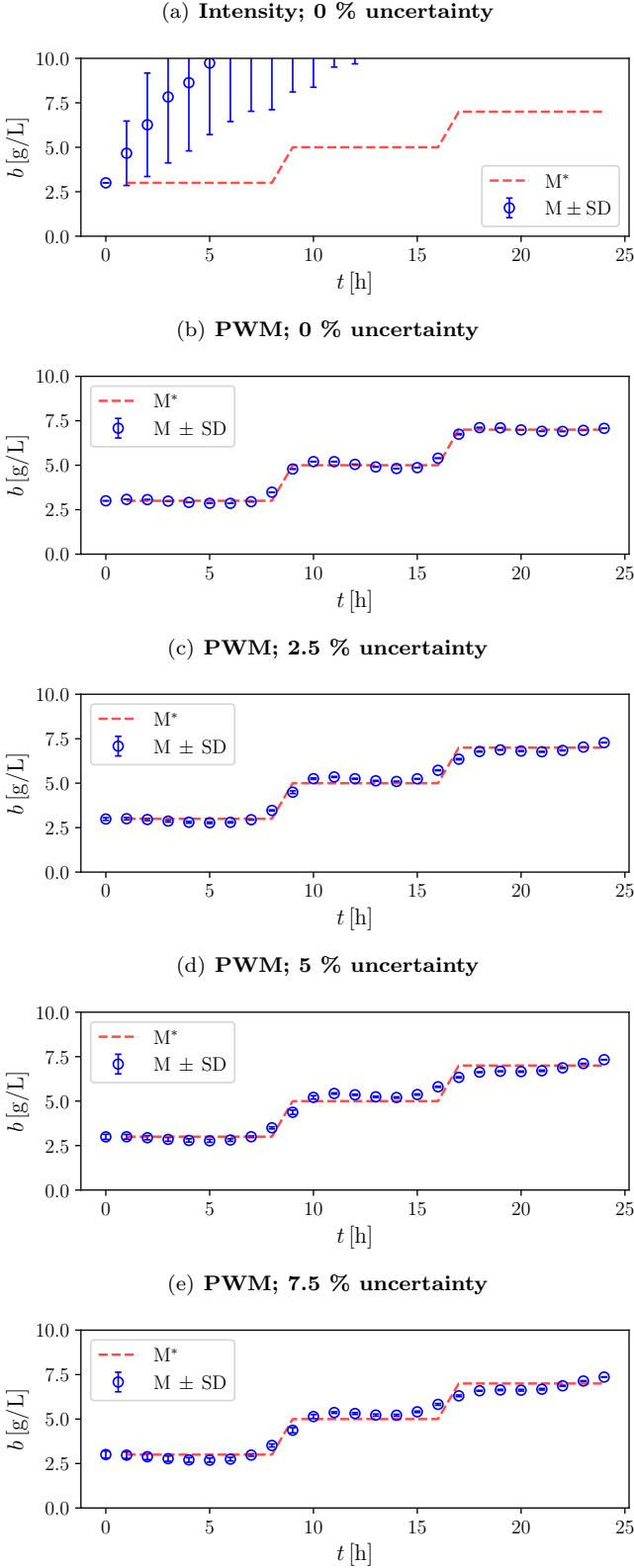
**(e) PWM; 7.5 % uncertainty**

Fig. 3. Optimized biomass trajectories versus the three-setpoint reference trajectory obtained with RL policies under intensity-driven and PWM-driven optogenetic control for selected uncertainty levels. M: mean; SD: standard deviation; M*: reference.



**(a) Intensity; 0 % uncertainty**

**(b) PWM; 0 % uncertainty**

**(c) PWM; 2.5 % uncertainty**

**(d) PWM; 5 % uncertainty**
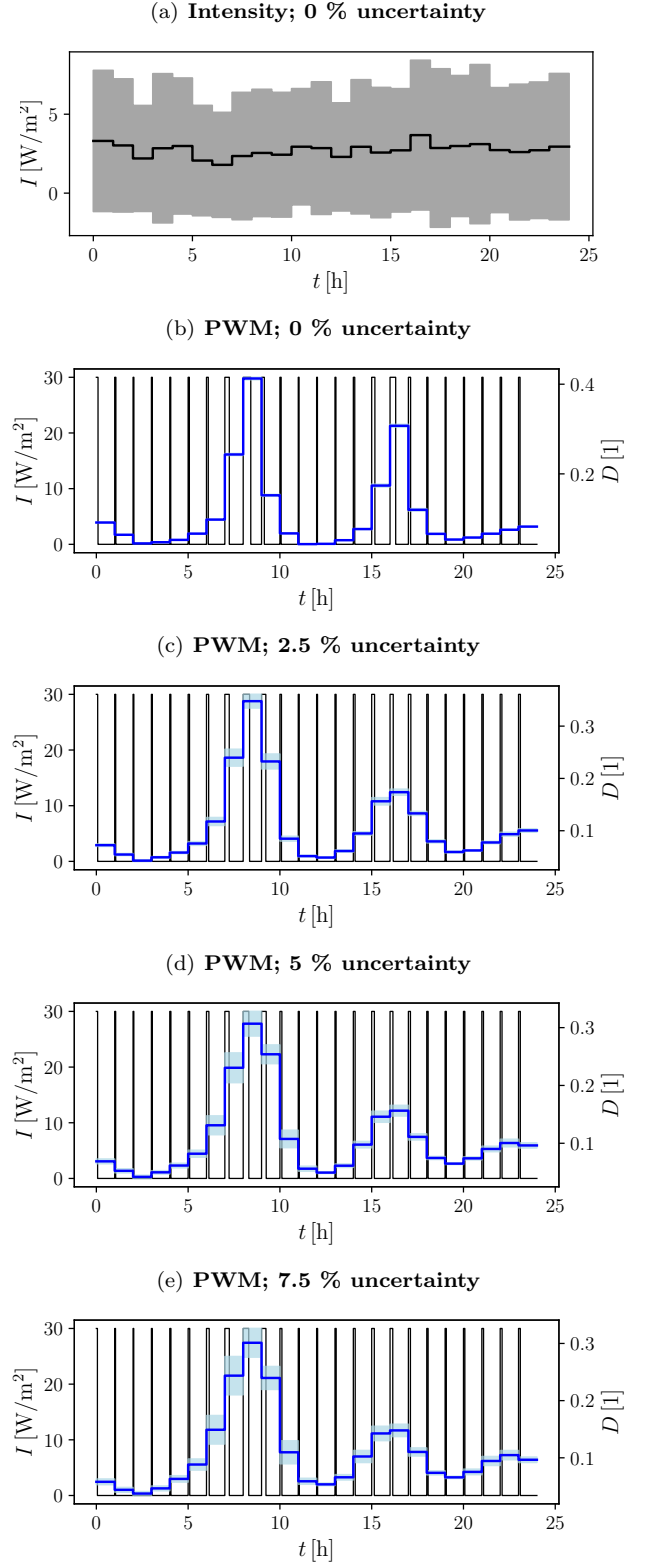
**(e) PWM; 7.5 % uncertainty**

Fig. 4. Optimized light input trajectories obtained with RL policies under intensity-driven and PWM-driven optogenetic control for selected uncertainty levels. The duty cycle mean is denoted by a blue line; the physical light intensity mean is denoted by a black line. In (a), the standard deviation of the intensity is represented as a shaded blue area. In (b)-(e), the standard deviation of the duty cycle is denoted as a shaded blue area.

ation does cover the range $[0, q_{\mathrm{p,max}}]$, despite the steepness. Note that we do not plot the intensity-driven benchmark

for the uncertainty scenarios, as it failed to converge even in the most optimistic case without uncertainty.

In contrast, the duty-cycle control, both without and with system uncertainty, led in all cases to high return values and clear convergence of the learning process (cf. Fig. 2 (b)-(e)), demonstrating the effectiveness and robustness of the proposed PWM-driven RL approach. The smoother tunability provided by PWM allows the policy to experience more stable and informative gradients during learning. The converged PWM-driven policies also exhibited low variability, which can be explained by the fact that, in essence, the duty cycle operates between only two modes (light ON/OFF within each forcing period). This makes the process less sensitive to uncertainties in intermediate gene activation levels.

The poor convergence of the intensity-driven optogenetic bioprocess is clearly reflected in Fig. 3, which compares the optimized biomass trajectories against the three-setpoint reference trajectory. In contrast, the PWM-driven bioprocess achieves successful tracking in all cases (with and without uncertainty), closely matching the reference. As expected, the PWM-driven process with no uncertainty yields practically perfect reference tracking, consistent with the return convergence behavior discussed above. Moreover, the low variability in the converged returns translates into very low variability in the biomass dynamics, even at the largest uncertainty level. As already mentioned, this is related to the fact that the process effectively operates between two robust modes (light ON/OFF within each forcing period), which are expected to be less sensitive to system uncertainty than intermediate gene-expression levels (as opposed to intensity-driven control). That said, as uncertainty increases, PWM-based control begins to struggle slightly at the setpoint transitions, yet overall tracking performance remains good.

Finally, taking a closer look at the optimized input trajectories (Fig. 4), we see that the intensity-driven control remains practically constant in terms of its mean intensity, but with a very large standard deviation, which is consistent with the poor learning convergence discussed above. In contrast, the PWM-driven control scenarios effectively modulate the duty cycles to match the reference trajectory. The duty cycles follow a coherent pattern across forcing periods, oscillating from lower to higher values to both maintain setpoints and enable setpoint transitions when appropriate. There is some increased standard deviation in the duty-cycle trajectories with rising uncertainty, yet it does not grow significantly. Overall, this demonstrates that our RL approach successfully learns continuous duty cycles that encode binary ON-OFF intensity patterns within forcing periods, resulting in robust PWM-driven control policies in the context of optogenetic bioprocesses.

## 5. CONCLUSION

In this study, we introduced an RL strategy to solve switching-time optimal control problems tailored to PWM optogenetics. The approach leverages duty cycles as continuous decision variables, which are decoded into alternating binary light-intensity inputs at the process level. On this basis, we use policy gradients to optimize the control policy. A case study involving optogenetic growth control in the presence of a steep light-gene-expression dose-response demonstrated the strong performance of the RL-derived PWM control policies, both in terms of tunability and robust tracking behavior under different levels of system uncertainty. Beyond optogenetics, the proposed control approach can be generalized to other bioprocesses involving duty-cycle-type inputs, such as microalgal photocycles or pulsed/intermittent feeding. Ongoing work focuses on extending the methodology to broader process contexts and to systems with a larger number of control inputs and degrees of freedom.

## REFERENCES

Benisch, M., Aoki, S.K., and Khammash, M. (2024). Unlocking the potential of optogenetics in microbial applications. *Current Opinion in Microbiology*, 77, 102404.

Benzinger, D., Ovinnikov, S., and Khammash, M. (2022). Synthetic gene networks recapitulate dynamic signal decoding and differential gene expression. *Cell Systems*, 13(5), 353–364.e6.

Davidson, E.A., Basu, A.S., and Bayer, T.S. (2013). Programming Microbes Using Pulse Width Modulation of Optical Signals. *Journal of Molecular Biology*, 425(22), 4161–4166.

Espinel-Ríos, S., Avalos, J.L., Del Rio Chanona, E.A., and Zhang, D. (2025a). Reinforcement learning for efficient and robust multi-setpoint and multi-trajectory tracking in bioprocesses. *Computers & Chemical Engineering*, 202, 109297.

Espinel-Ríos, S., Walser, R., and Zhang, D. (2025b). Reinforcement Learning for Robust Dynamic Metabolic Control. *Biotechnology and Bioengineering*, bit.70077.

Ewing, T.A., Nouse, N., Van Lint, M., Van Haveren, J., Hugenholtz, J., and Van Es, D.S. (2022). Fermentation for the production of biobased chemicals in a circular economy: a perspective for the period 2022–2050. *Green Chemistry*, 24(17), 6373–6405.

Hoffman, S.M., Tang, A.Y., and Avalos, J.L. (2022). Optogenetics illuminates applications in microbial engineering. *Annual Review of Chemical and Biomolecular Engineering*, 13(1), 373–403.

Konzock, O. and Nielsen, J. (2024). TRYing to evaluate production costs in microbial biotechnology. *Trends in Biotechnology*, 42(11), 1339–1347.

Milias-Argeitis, A., Rullan, M., Aoki, S.K., Buchmann, P., and Khammash, M. (2016). Automated optogenetic feedback control for precise and robust regulation of gene expression and cell growth. *Nature Communications*, 7(1), 12546.

Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., and Chintala, S. (2019). *PyTorch: an imperative style, high-performance deep learning library*. Curran Associates Inc., Red Hook, NY, USA.

Sutton, R.S., McAllester, D., Singh, S., and Mansour, Y. (1999). Policy gradient methods for reinforcement learning with function approximation. In S. Solla, T. Leen, and K. Müller (eds.), *Advances in Neural Information Processing Systems*, volume 12. MIT Press.