

Hybrid Model Assisted Reinforcement Learning for Cell Therapy Manufacturing Process Control

Hua Zheng¹ Wei Xie¹ Keqi Wang¹ Zheng Li²

¹Department of Mechanical and Industrial Engineering
Northeastern University,

²Genentech, Inc., South San Francisco, CA, USA

Challenges in Cell Therapy manufacturing

- **High complexity:**

Challenges in Cell Therapy manufacturing

- **High complexity:**
 - The productivity and functional identity of cell products are sensitive to cell culture conditions.

Challenges in Cell Therapy manufacturing

- **High complexity:**

- The productivity and functional identity of cell products are sensitive to cell culture conditions.
- Improper cultivation can not only hinder yield, but can result in heterogeneously differentiated cell populations.

Challenges in Cell Therapy manufacturing

- **High complexity:**

- The productivity and functional identity of cell products are sensitive to cell culture conditions.
- Improper cultivation can not only hinder yield, but can result in heterogeneously differentiated cell populations.

- **Limited data:**

Challenges in Cell Therapy manufacturing

- **High complexity:**

- The productivity and functional identity of cell products are sensitive to cell culture conditions.
- Improper cultivation can not only hinder yield, but can result in heterogeneously differentiated cell populations.

- **Limited data:**

- ① Lengthy analytical testing time for complex cell therapeutics

Challenges in Cell Therapy manufacturing

- **High complexity:**

- The productivity and functional identity of cell products are sensitive to cell culture conditions.
- Improper cultivation can not only hinder yield, but can result in heterogeneously differentiated cell populations.

- **Limited data:**

- ① Lengthy analytical testing time for complex cell therapeutics
- ② More and more personalized cell therapeutics

Challenges in Cell Therapy manufacturing

- **High complexity:**

- The productivity and functional identity of cell products are sensitive to cell culture conditions.
- Improper cultivation can not only hinder yield, but can result in heterogeneously differentiated cell populations.

- **Limited data:**

- ① Lengthy analytical testing time for complex cell therapeutics
- ② More and more personalized cell therapeutics

- **High variability:** seed cells can be extracted and isolated from individual patients and donors, which leads to high variability

Limitation of State-of-the-art Methods

Existing mechanistic models often ignore various sources of process stochasticity:

- batch-to-batch variation (Mockus et al., 2015)
- intracellular production fluctuations (Vasdekis et al., 2015)
- Raw material variability (Dickens et al., 2018).

Limitation of State-of-the-art Methods

Existing mechanistic models often ignore various sources of process stochasticity:

- batch-to-batch variation (Mockus et al., 2015)
- intracellular production fluctuations (Vasdekis et al., 2015)
- Raw material variability (Dickens et al., 2018).

For classical control and reinforcement learning control

- classical control strategies are often derived from deterministic mechanistic models and overlook bioprocess stochastic uncertainty & model uncertainty
- RL approaches often do not have good way to incorporate enough prior knowledge on bioprocessing mechanisms;

Hybrid-RL with Probabilistic Knowledge Graph

Driven by the critical challenges, we propose a data-driven stochastic optimization framework named “hybrid-RL”.

- **KG network hybrid model** is probabilistic and mechanism-based and created to characterize the spatial-temporal causal interdependencies between critical process parameters (CPPs) and critical quality attributes (CQAs).
- **Bayesian inference** is used to derive a posterior distribution of the hybrid model.
- **Hybrid model-based Bayesian RL** (called “**hybrid-RL**”) is developed to efficiently guide optimal, robust, and interpretable dynamic decision making.

The proposed Hybrid-RL framework demonstrates promising performance for cell therapy manufacturing process optimization.

Problem Statement and Bayesian RL

We model the cell therapy manufacturing process as a finite-horizon Markov decision process (MDP) specified by $(\mathcal{S}, \mathcal{A}, H, r, p)$.

The system start at an initial state \mathbf{s}_1 drawn from $p_1(\mathbf{s}_1)$. At any time t ,

- the agent observes the state $\mathbf{s}_t \in \mathcal{S}$ and takes an action $\mathbf{a}_t \in \mathcal{A}$ from a policy $\pi_t(\mathbf{s}_t | \mathbf{a}_t)$.
- receives a reward $r_t(\mathbf{s}_t, \mathbf{a}_t) \in \mathbb{R}$.

Thus, the probabilistic model of the process trajectory $\boldsymbol{\tau} = (\mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_H, \mathbf{a}_H, \mathbf{s}_{H+1})$, i.e.,

$$p(\boldsymbol{\tau} | \boldsymbol{\theta}) = p(\mathbf{s}_1) \prod_{t=1}^H p(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t; \boldsymbol{\theta}_t),$$

Given $\boldsymbol{\theta}$, the **performance** of the policy π is evaluated via the expected accumulated reward,

$$J(\pi; \boldsymbol{\theta}) \equiv \mathbb{E}_{\boldsymbol{\tau}} \left[\sum_{t=1}^{H+1} r_t(\mathbf{s}_t, \mathbf{a}_t) \mid \pi, \boldsymbol{\theta} \right], \quad (1)$$

Problem Statement (Cont'd)

Latent state: Let \mathbf{z}_t denote the latent state variables. Thus, at any time step t ,

- We have observable and unobservable state $\mathbf{s}_t = (\mathbf{x}_t, \mathbf{z}_t)$.

- We have the likelihood of the partially observed trajectory

$$\boldsymbol{\tau}_x \equiv (\mathbf{x}_1, \mathbf{a}_1, \dots, \mathbf{x}_H, \mathbf{a}_H, \mathbf{x}_{H+1}) \text{ is } p(\boldsymbol{\tau}_x | \boldsymbol{\theta}) = \int \cdots \int p(\boldsymbol{\tau} | \boldsymbol{\theta}) d\mathbf{z}_1 \cdots d\mathbf{z}_{H+1}.$$

Model uncertainty is quantified by a posterior distribution obtained by applying Bayesian rule,

$$p(\boldsymbol{\theta} | \mathcal{D}) \propto p(\boldsymbol{\theta}) P(\mathcal{D} | \boldsymbol{\theta}) = p(\boldsymbol{\theta}) \prod_{i=1}^m p(\boldsymbol{\tau}_x^{(i)} | \boldsymbol{\theta}) \quad (2)$$

where the prior $p(\boldsymbol{\theta})$ can incorporate the mechanism knowledge on the model parameters.

Objective: the optimization problem of KG hybrid model-based Bayesian RL is formulated by

$$\pi^* = \arg \max_{\pi \in \mathcal{P}} \mathcal{J}(\pi) \quad (3)$$

with \mathcal{P} representing the feasible set of decision policies and the optimization objective

$$\mathcal{J}(\pi) \equiv \mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta} | \mathcal{D})} [J(\pi; \boldsymbol{\theta})]$$

with (1) inner expectation in $J(\pi; \boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\tau}} [\sum_{t=1}^{H+1} r_t(\mathbf{s}_t, \mathbf{a}_t) | \pi, \boldsymbol{\theta}]$ accounting for inherent stochasticity; and (2) outer expectation accounting for model uncertainty.

Bioprocess Hybrid Modeling

- Given the existing ODE-based mechanistic model $ds/dt = \mathbf{f}(\mathbf{s}, \mathbf{a}; \boldsymbol{\beta})$, we construct the hybrid model for state transition,

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \Delta t \cdot \mathbf{f}_x(\mathbf{x}_t, \mathbf{z}_t, \mathbf{a}_t; \boldsymbol{\beta}_t) + \mathbf{e}_{t+1}^x, \quad (4)$$

$$\mathbf{z}_{t+1} = \mathbf{z}_t + \Delta t \cdot \mathbf{f}_z(\mathbf{x}_t, \mathbf{z}_t, \mathbf{a}_t; \boldsymbol{\beta}_t) + \mathbf{e}_{t+1}^z, \quad (5)$$

where the residual terms $\mathbf{e}_{t+1}^x \sim \mathcal{N}(0, V_{t+1}^x)$ and $\mathbf{e}_{t+1}^z \sim \mathcal{N}(0, V_{t+1}^z)$.

- Let $g(\mathbf{s}_t, \mathbf{a}_t; \boldsymbol{\beta}_t) \equiv \mathbf{s}_t + \Delta t \cdot \mathbf{f}(\mathbf{s}_t, \mathbf{a}_t; \boldsymbol{\beta}_t)$. At any time step $t \in \mathcal{H}$, we have

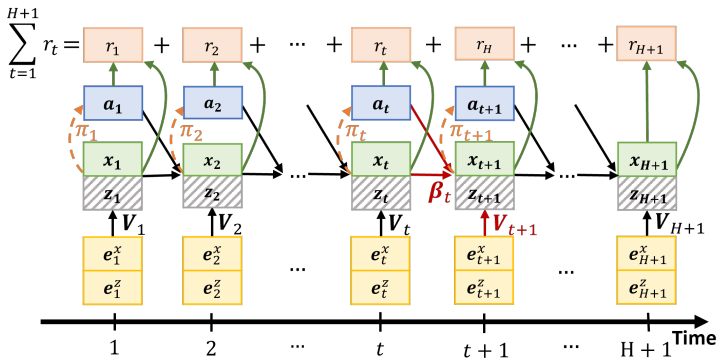
$$\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t = g(\mathbf{s}_t, \mathbf{a}_t; \boldsymbol{\beta}_t) + \mathbf{e}_{t+1} \sim \mathcal{N}\left(g(\mathbf{s}_t, \mathbf{a}_t; \boldsymbol{\beta}_t), V_{t+1}\right) \quad (6)$$

where $\mathbf{s}_t = (\mathbf{x}_t, \mathbf{z}_t)$, $\mathbf{e}_{t+1} = (\mathbf{e}_{t+1}^x, \mathbf{e}_{t+1}^z)$, and V_{t+1} is diagonal covariance matrix with diagonal entries from V_{t+1}^x and V_{t+1}^z .

- Random mechanistic coefficients $\boldsymbol{\beta}_t$ account for batch-to-batch variation.
 $\boldsymbol{\theta}_t = (\boldsymbol{\mu}_t^\beta, \text{vec}(\boldsymbol{\Sigma}_t^\beta), \text{vec}(V_{t+1}^x), \text{vec}(V_{t+1}^z))^\top$.

Bioprocess Hybrid Modeling (Cont'd)

A Second View of Hybrid-RL: Policy-Augmented Knowledge Graph (KG)



Bayesian Inference

Inference: ABC-SMC sampling procedure for generating posterior samples from $p(\boldsymbol{\theta}|\mathcal{D})$ (derived from Toni et al. (2009); Lenormand et al. (2013); Del Moral et al. (2006)).

Main Idea: Given the observed trajectory $\boldsymbol{\tau}_x = (\mathbf{x}_1, \mathbf{a}_1, \dots, \mathbf{x}_H, \mathbf{a}_H, \mathbf{x}_{H+1})$.

$$p(\boldsymbol{\theta}|\boldsymbol{\tau}_x) \propto p(\boldsymbol{\tau}_x|\boldsymbol{\theta})p(\boldsymbol{\theta}) \quad (7)$$

The algorithm samples $\boldsymbol{\theta}$ and $\boldsymbol{\tau}_x^*$ from the joint posterior:

$$p_\delta(\boldsymbol{\theta}, \boldsymbol{\tau}_x^*|\boldsymbol{\tau}) = \frac{p(\boldsymbol{\theta})p(\boldsymbol{\tau}_x|\boldsymbol{\theta})\mathbb{1}_\delta[\boldsymbol{\tau}_x^*]}{\int \int p(\boldsymbol{\theta})p(\boldsymbol{\tau}_x^*|\boldsymbol{\theta})\mathbb{1}_\delta[\boldsymbol{\tau}_x^*]d\boldsymbol{\tau}_x^*d\boldsymbol{\theta}} \quad (8)$$

where $\mathbb{1}_\delta[\boldsymbol{\tau}_x^*] = \mathbb{1}_\delta[d(\boldsymbol{\tau}_x, \boldsymbol{\tau}_x^*) \leq \delta]$ is one if $d(\boldsymbol{\tau}_x, \boldsymbol{\tau}_x^*) \leq \delta$ and zero else.

When δ is small, $p_\delta(\boldsymbol{\theta}|\boldsymbol{\tau}_x) = \int p_\delta(\boldsymbol{\theta}, \boldsymbol{\tau}_x^*|\boldsymbol{\tau})d\boldsymbol{\tau}_x^*$ is a good approximation to $p(\boldsymbol{\theta}|\boldsymbol{\tau}_x)$

KG Hybrid Model-based Bayesian RL

For each $t \in \mathcal{H}$, we define the state value function $V_t^\pi(\mathbf{s}) : \mathcal{S} \rightarrow \mathbb{R}$ and action value function $Q_t^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ as

$$V_t^\pi(\mathbf{s}) = \mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D})} \mathbb{E}_{p(\mathbf{s}_{t+1}|\mathbf{s}_t, \pi_t(\mathbf{s}_t); \boldsymbol{\theta}_t)} \left[\sum_{\ell=t}^H r_\ell(\mathbf{s}_\ell, \pi_\ell(\mathbf{s}_\ell)) \mid \mathbf{s}_t = \mathbf{s} \right]$$
$$Q_t^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E} \left[\sum_{\ell=t}^H r_\ell(\mathbf{s}_\ell, \pi_\ell(\mathbf{s}_\ell)) \mid \mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a} \right]$$
$$= r_t(\mathbf{s}, \mathbf{a}) + \mathbb{E} [V_{t+1}^\pi(\mathbf{s}_{t+1}) \mid \mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a}] \quad (9)$$

The Bellman optimality equation (Sutton and Barto, 2018, Chapter 3.6)

$$V_t^*(\mathbf{s}) = \max_{\mathbf{a} \in \mathcal{A}} r_t(\mathbf{s}, \mathbf{a}) + \mathbb{E} [V_{t+1}^*(\mathbf{s}_{t+1}) \mid \mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a}]$$
$$= \max_{\mathbf{a} \in \mathcal{A}} Q_t^*(\mathbf{s}, \mathbf{a}). \quad (10)$$

The optimal greedy policy (Puterman, 2014) with

$$\pi_t^*(\mathbf{s}) \equiv \operatorname{argmax}_{\mathbf{a} \in \mathcal{A}} Q_t^*(\mathbf{s}, \mathbf{a}), \text{ for any } \mathbf{s} \in \mathcal{S}. \quad (11)$$

Bayesian Sparse Sampling (Kearns et al., 2002; Wang et al., 2005)

Input: state \mathbf{s}_t ; scenario numbers B and J for estimating $\mathbb{E}_{p(\boldsymbol{\theta}|\mathcal{D})} \mathbb{E}_{p(\mathbf{s}_{t+1}|\mathbf{s}_t, \pi(\mathbf{s}_t; \boldsymbol{\theta}_t))}[\cdot]$; $\widehat{p}(\boldsymbol{\theta}|\mathcal{D})$ from SMC-ABC.

Output: Estimated optimal Q-function $\widehat{Q}(\mathbf{s}, \mathbf{a})$

Function $Q_{\text{FUN}}(t, \mathbf{s}_t, \mathbf{a}_t)$:

for $b = 1, 2, \dots, B$ do

(A1) Generate a posterior sample of model parameters, $\boldsymbol{\theta}_b \sim \widehat{p}(\boldsymbol{\theta}_t|\mathcal{D})$.

 for $j = 1, \dots, J$ do

(A2) Sample from state transition distribution, $\mathbf{s}_{t+1}^{(b,j)} \sim p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t; \boldsymbol{\theta}_t, b)$

(A3) $V_{t+1}(\mathbf{s}_{t+1}^{(b,j)}) = V_{\text{FUN}}(t+1, \mathbf{s}_{t+1}^{(b,j)})$

(A4) $\widehat{Q}_t(\mathbf{s}_t, \mathbf{a}_t) = r_t(\mathbf{s}_t, \mathbf{a}_t) + \frac{1}{BJ} \sum_{b=1}^B \sum_{j=1}^J V_{t+1}(\mathbf{s}_{t+1}^{(b,j)})$.

return $\widehat{Q}_t(\mathbf{s}_t, \mathbf{a}_t)$.

Function $V_{\text{FUN}}(t, \mathbf{s}_t)$:

if $t = H + 1$ then

 return $r_{H+1}(\mathbf{s}_{H+1})$;

for $\mathbf{a}_t \in \mathcal{A}$ do

 for $b = 1, 2, \dots, B$ do

(B1) Generate a posterior sample of model parameters $\boldsymbol{\theta}_b \sim \widehat{p}(\boldsymbol{\theta}_t|\mathcal{D})$

 for $j = 1, \dots, J$ do

(B2) Sample from state transition $\mathbf{s}_{t+1}^{(b,j)} \sim p(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t; \boldsymbol{\theta}_t, b)$

(B3) $V_{t+1}(\mathbf{s}_{t+1}^{(b,j)}) = V_{\text{FUN}}(t+1, \mathbf{s}_{t+1}^{(b,j)})$

(B4) Estimate $\widehat{Q}_t(\mathbf{s}_t, \mathbf{a}_t) = r_t(\mathbf{s}_t, \mathbf{a}_t) + \frac{1}{BJ} \sum_{b=1}^B \sum_{j=1}^J V_{t+1}(\mathbf{s}_{t+1}^{(b,j)})$

(B5) $\widehat{V}_t(\mathbf{s}_t) = \max_{\mathbf{a}_t \in \mathcal{A}} \widehat{Q}_t(\mathbf{s}_t, \mathbf{a}_t)$ as in (10)

Cell Therapy Manufacturing Case Study

Mechanistic Model: Glen et al. (2018) developed an ODE-based mechanistic model describing the dynamics of an unidentified autocrine growth inhibitor accumulation and its impact on the erythroblast cell production process. We extend this model to a two phases: growth and stationary phases with index $p = 1, 2$.

$$\frac{d\rho}{dt} = r_p^g \rho \left(1 - \left(1 + e^{(k_p^s(k_p^c - I))} \right)^{-1} \right), \quad (12)$$

$$\frac{dI}{dt} = \frac{d\rho}{dt} - r_p^d I, \quad (13)$$

- ρ_t and I_t represent the cell density and the inhibitor concentration at t .
- The kinetic coefficients r_p^g , k_p^s , k_p^c and r_p^d denote the cell growth rate, inhibitor sensitivity, inhibitor threshold, and inhibitor decay. The phase transition occurs at $T_\star = 18$ hour.

Simulator: based on (12)-(13), we develop a simulator by including various source of uncertainty

$$d\rho = r_p^g \rho \left(1 - \left(1 + e^{(k_p^s(k_p^c - I))} \right)^{-1} \right) dt + \sigma_n dW \quad (14)$$

$$dI = d\rho - r_p^d I dt + \sigma_n dW \quad (15)$$

with (1) **random initial values** $\rho_1 \sim \mathcal{N}(\mu_\rho, \sigma_\rho^2)$ and $I_1 = 0$, (2) **batch-to-batch variation** $r_p^g \sim \mathcal{N}(\mu_p^g, (\sigma_p^g)^2)$ with $p = 1, 2$, and **measurement error** $\rho_t \leftarrow \rho_t + e_m$, with $e_m \sim \mathcal{N}(0, \sigma_m^2)$.

Prediction Error (MAE) of Cell Density

Assessment of the long-term prediction performance (mean absolute error) of the KG hybrid model and the ODE mechanistic model fitted by LS method (LS-ODE).

Both models were fitted by “real-world” historical trajectories with the size $m = 3, 6, 20$. We evaluate performance based on $r = 30$ macro-replications.

Noise Level		h (hrs)	Hybrid			LS-ODE		
b2b	noise		$m = 3$	$m = 6$	$m = 20$	$m = 3$	$m = 6$	$m = 20$
high	$\sigma_n = 0.01$	3	0.12 \pm 0.05	0.09 \pm 0.03	0.06 \pm 0.02	0.41 \pm 0.19	0.59 \pm 0.30	0.44 \pm 0.22
		18	0.60 \pm 0.17	0.48 \pm 0.10	0.26 \pm 0.07	0.74 \pm 0.15	0.57 \pm 0.25	0.49 \pm 0.23
		30	0.59 \pm 0.16	0.40 \pm 0.11	0.22 \pm 0.06	0.65 \pm 0.24	0.70 \pm 0.36	0.84 \pm 0.65
high	$\sigma_n = 0.03$	3	0.21 \pm 0.05	0.14 \pm 0.03	0.08 \pm 0.03	0.37 \pm 0.20	0.40 \pm 0.19	0.36 \pm 0.23
		18	1.07 \pm 0.22	0.82 \pm 0.15	0.48 \pm 0.12	1.11 \pm 0.28	0.93 \pm 0.31	0.83 \pm 0.34
		30	1.11 \pm 0.24	0.74 \pm 0.16	0.44 \pm 0.12	1.57 \pm 0.76	1.09 \pm 0.41	0.93 \pm 0.45
low	$\sigma_n = 0.01$	3	0.10 \pm 0.03	0.07 \pm 0.02	0.04 \pm 0.01	0.38 \pm 0.23	0.54 \pm 0.26	0.38 \pm 0.20
		18	0.48 \pm 0.12	0.38 \pm 0.09	0.27 \pm 0.06	0.43 \pm 0.13	0.35 \pm 0.12	0.28 \pm 0.11
		30	0.47 \pm 0.11	0.30 \pm 0.08	0.16 \pm 0.04	0.45 \pm 0.11	0.52 \pm 0.23	0.32 \pm 0.15
low	$\sigma_n = 0.03$	3	0.18 \pm 0.06	0.13 \pm 0.03	0.04 \pm 0.01	0.68 \pm 0.34	0.46 \pm 0.23	0.31 \pm 0.18
		18	1.00 \pm 0.20	0.69 \pm 0.15	0.27 \pm 0.06	1.27 \pm 0.36	1.47 \pm 1.41	0.55 \pm 0.16
		30	1.04 \pm 0.28	0.65 \pm 0.17	0.16 \pm 0.04	1.40 \pm 0.48	1.61 \pm 1.40	0.66 \pm 0.20

Remark: b2b is batch-to-batch variation and noise is the process noise.

Medium Full Exchange Decision Making

Background: Medium exchange is an essential element of successful long-term cell culture. Culture medium is exchanged to supply new nutrients and to eliminate waste products produced by the cells.

Goal: finding the optimal time to fully exchange the medium with fresh medium.

State: is defined as cell density and inhibitor $\mathbf{s}_t = (\rho_t, I_t)$.

Action: $a_t = 0$ denoting the full exchange of medium at step t ; $a_t = 1$, otherwise. $I_t = a_t I_t$ represent the post-exchanged concentration of inhibitor

Reward:

- Total operational cost: $C(T, M) = C_t T + C_m M$,
- Reward function is defined by cell yield per cost — the efficiency of the system during the T hours (H time steps) cell culture (Glen et al., 2018):

$$r_t = 0 \quad \text{with} \quad 0 \leq t \leq H$$

$$r_{H+1}(\mathbf{s}_{H+1}, a_{H+1} = \text{"Harvest"}) = \frac{M(\rho_T - \rho_0)}{C(T, M)}$$

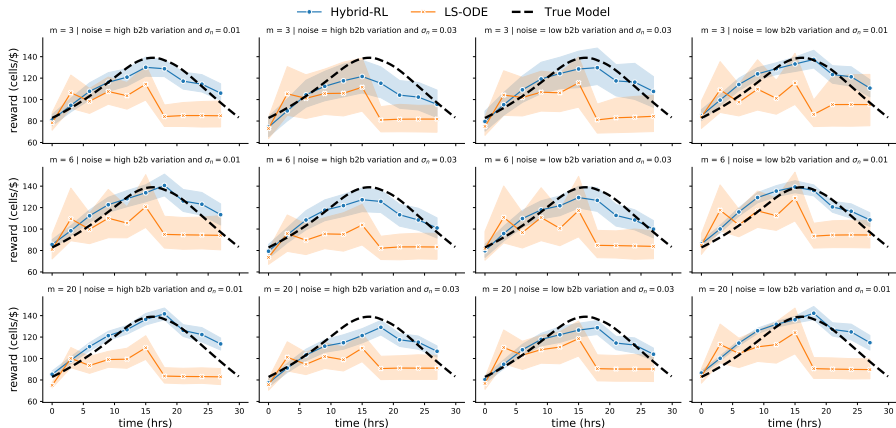
where ρ_T represents the cell density at the T -th hour.

Decision Epochs: The medium is fully exchanged up to one time in decision hours $\{0, 3, \dots, 27\}$.

Medium Exchange Cost Efficiency

Performance of hybrid-RL and LS-ODE in 30 macro-replications.

- The validated models are used to optimize the media exchange time for cells to be produced with optimal cost efficiency at a given production scale (100L).
- The number of cells produced per dollar (y-axis) for a given time point (x-axis) of media exchange are calculated for \$150/hr operating time cost and \$10/L of consumable cost



Cell Culture Expansion Scheduling

- Each expansion, the original batch is scaled up to a n times larger cell culture vessel filling with fresh medium.
- Cell density ρ and the concentration of inhibitor I decrease to $1/n$ of original batch immediately after each scale-up.
- The reward function is then defined by the difference between revenue and cost as,

$$r_t = 0 \text{ with } 0 \leq t \leq H$$

$$r_{H+1}(\mathbf{s}_{H+1}, \mathbf{a}_{H+1} = \text{"Harvest"})$$

$$= K(\rho_T, \xi, n) - C(T, M)$$

where the revenue $K(\rho_T, \xi, n) = P_c \times \rho_T \times n^\xi$

Noise Level		Hybrid-RL			LS-ODE		
b2b	noise	$m = 3$	$m = 6$	$m = 20$	$m = 3$	$m = 6$	$m = 20$
high	$\sigma_n = 0.01$	7317.24 (352.37)	7588.15 (322.53)	7892.84 (69.65)	5677.62 (389.24)	5944.25 (403.52)	6030.83 (399.61)
high	$\sigma_n = 0.03$	6888.86 (693.07)	7266.23 (393.66)	7689.45 (143.11)	-2259.01 (704.82)	1026.60 (484.50)	2454.22 (262.15)
low	$\sigma_n = 0.01$	7800.60 (151.01)	7955.38 (75.78)	8035.71 (52.48)	6115.31 (413.23)	6193.70 (381.83)	6417.39 (389.06)
low	$\sigma_n = 0.03$	7414.34 (334.41)	7572.99 (329.62)	7974.35 (76.15)	5978.52 (389.96)	6126.60 (393.61)	6225.02 (395.50)

Thank you!

- Pierre Del Moral, Arnaud Doucet, and Ajay Jasra. Sequential monte carlo samplers. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(3):411–436, 2006.
- Jason Dickens, Sarwat Khattak, Thomas E Matthews, Dave Kolwyck, and Kelly Wiltberger. Biopharmaceutical raw material variation and control. *Current opinion in chemical engineering*, 22:236–243, 2018.
- Katie E Glen, Elizabeth A Cheeseman, Adrian J Stacey, and Robert J Thomas. A mechanistic model of erythroblast growth inhibition providing a framework for optimisation of cell therapy manufacturing. *Biochemical Engineering Journal*, 133: 28–38, 2018.
- Michael Kearns, Yishay Mansour, and Andrew Y Ng. A sparse sampling algorithm for near-optimal planning in large markov decision processes. *Machine learning*, 49(2):193–208, 2002.
- Maxime Lenormand, Franck Jabot, and Guillaume Deffuant. Adaptive approximate bayesian computation for complex models. *Computational Statistics*, 28(6):2777–2796, 2013.
- Linas Mockus, John J Peterson, Jose Miguel Lainez, and Gintaras V Reklaitis. Batch-to-batch variation: a key component for modeling chemical manufacturing processes. *Organic Process Research & Development*, 19(8):908–914, 2015.
- Martin L Puterman. *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons, 2014.
- Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. A Bradford Book, Cambridge, MA, USA, 2018. ISBN 0262039249.
- Tina Toni, David Welch, Natalja Strelkowa, Andreas Ipsen, and Michael PH Stumpf. Approximate bayesian computation scheme for parameter inference and model selection in dynamical systems. *Journal of the Royal Society Interface*, 6(31): 187–202, 2009.
- AE Vasdekis, Andrew M Silverman, and Gregory Stephanopoulos. Origins of cell-to-cell bioprocessing diversity and implications of the extracellular environment revealed at the single-cell level. *Scientific Reports*, 5(1):1–7, 2015.
- Tao Wang, Daniel Lizotte, Michael Bowling, and Dale Schuurmans. Bayesian sparse sampling for on-line reward optimization. In *Proceedings of the 22nd international conference on Machine learning*, pages 956–963, 2005.