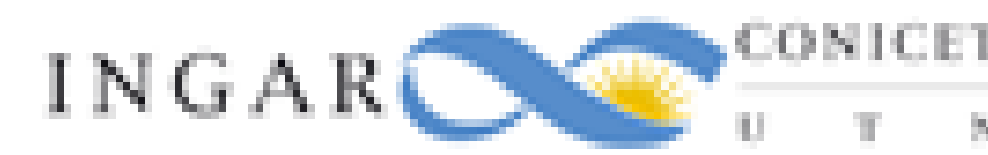


Probabilistic Modeling for Optimization of Bioreactors using Reinforcement Learning with Active Inference

Ernesto C. Martínez, Jong Woo Kim, Tilman Barz, M. Nicolás Cruz B.



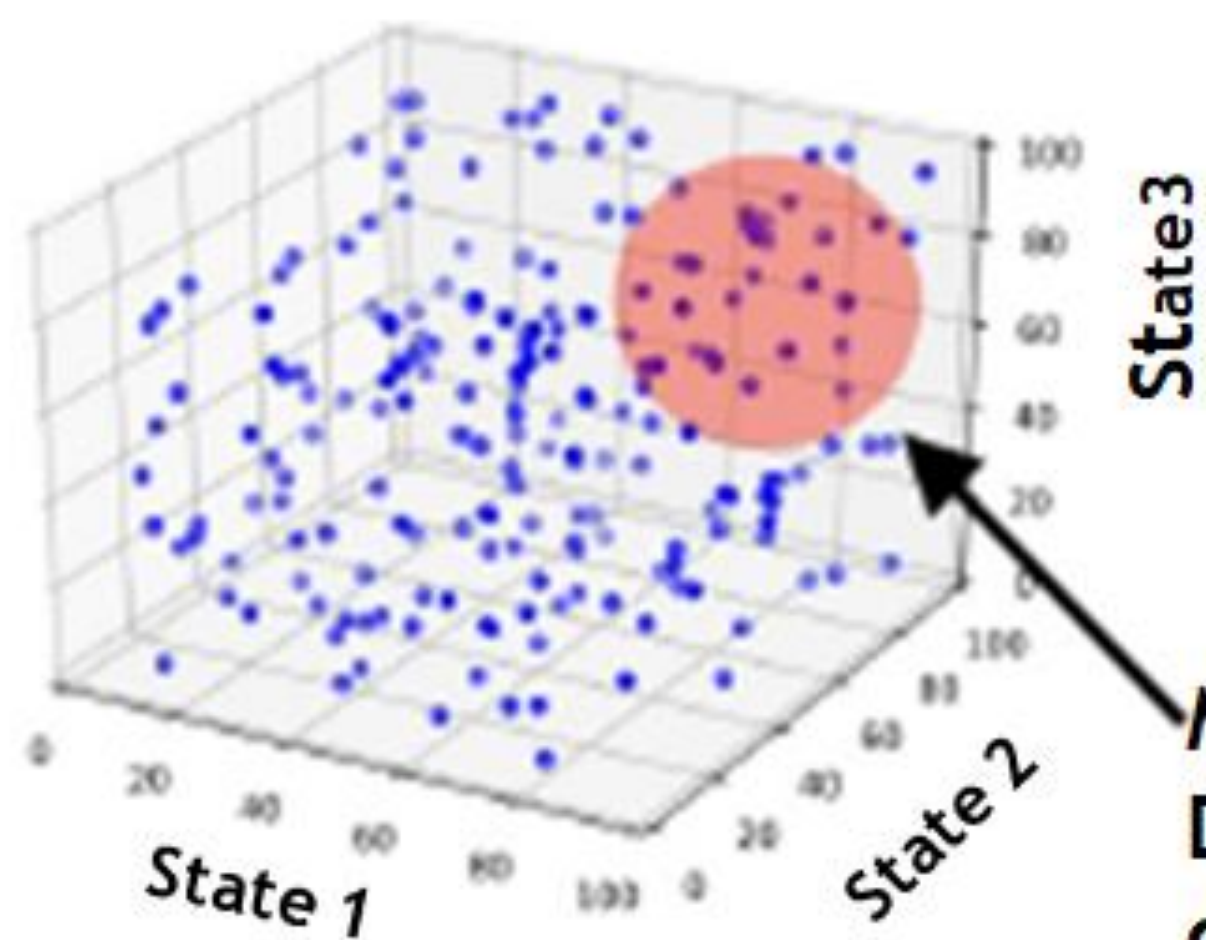
Technische Universität Berlin, Chair of Bioprocess Engineering, KIWI-biolab
www.bioprocess.tu-berlin.de
www.kiwi-biolab.de



Motivation and Scope

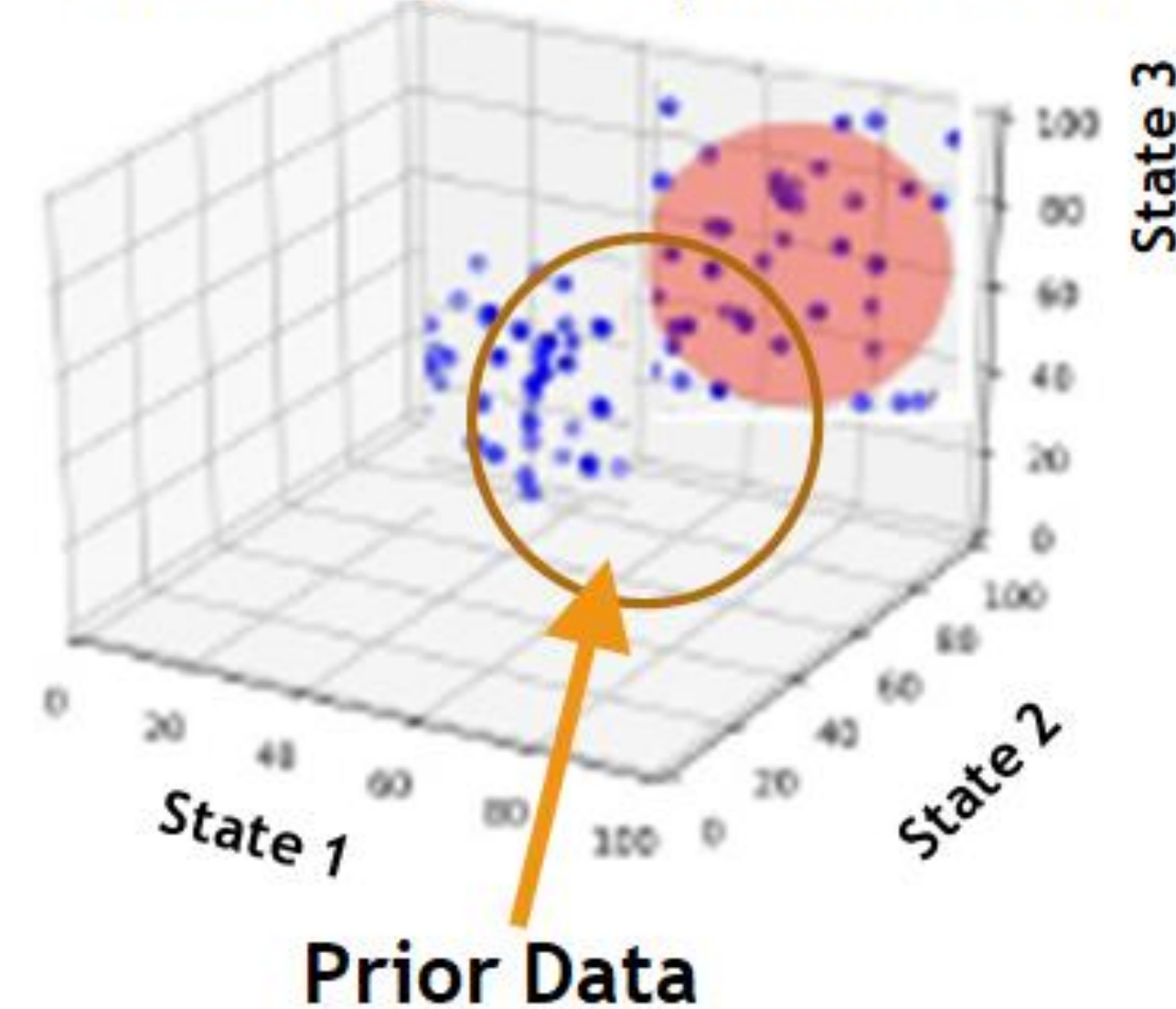
With **low-complexity imperfect** models, bioreactor optimization is often **sub-optimal**. Why?

Parametric precision



Even though is not widely accepted nor recognized, abstract (e. g., macroscopic or cybernetic) models used for bioreactor optimization are too shallow to account for the complexity of switching in metabolic pathways when a microorganism is responding to changes in the abiotic conditions

Modeling for optimization



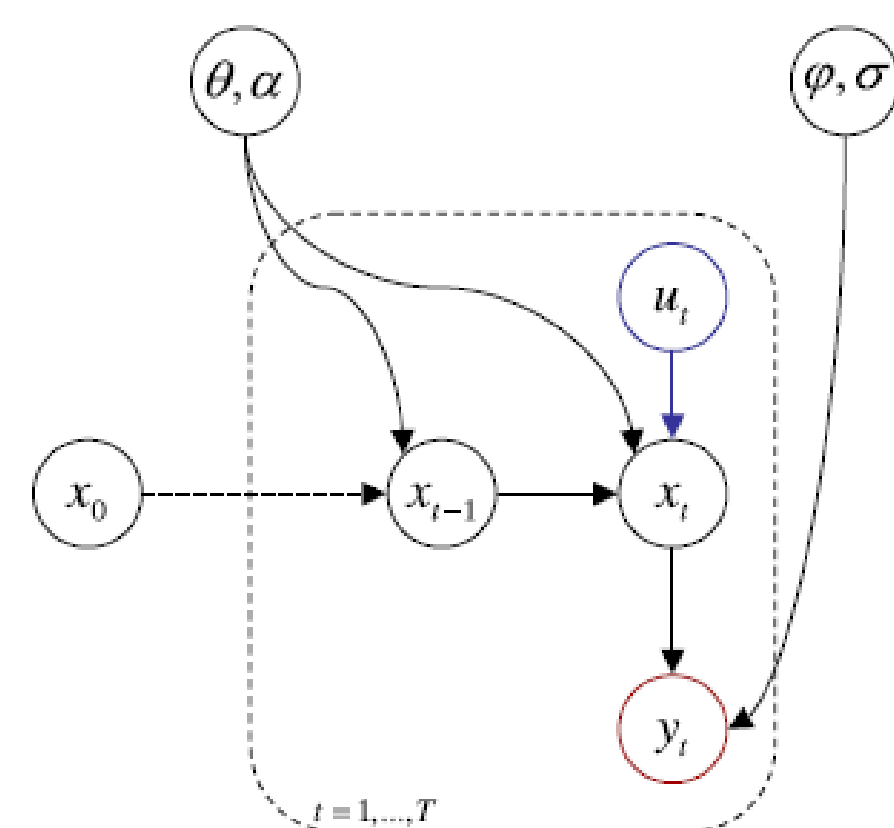
A challenge in *modeling for optimization* is how effective bioreactor models can be learned from designed experiments given (i) the rich complexity of profiling operating conditions, and (ii) the circular dependence of model learning and information content of sampled data, which often leads to suboptimal performance and low reproducibility.

Methodology and Results

Probabilistic (causal) models of bioreactors

A probabilistic (causal) model of a bioreactor is defined by a joint probability distribution over the following set of stochastic variables:

- x ; y : the $n \times n$, hidden states time-series; the $p \times n$, observations (sampled data),
- u : the $n_u \times n$, manipulated (controlled) inputs time-series,
- θ, φ : the $n_\theta \times 1$ evolution parameters; the $n_\varphi \times 1$ observation parameters,
- α : the state noise precision (structural errors),
- σ : the measurement noise precision (analytical and sensor calibration errors).



These variables are assumed to follow the (hidden) **state evolution and observation equations**:

$$x_t = f(x_{t-1}, \theta, u_{t-1}) + \eta_t; \quad \eta_t = N(0, \alpha^{-1}I)$$

$$y_t = g(x_t, \varphi) + \varepsilon_t; \quad \varepsilon_t = N(0, \sigma^{-1}I)$$

The **probabilistic model** m of a bioreactor is completely specified by the (initial) *Gaussian* prior distributions for its parameters θ, φ , and the *Gamma* priors for the precision hyper-parameters α, σ .

Adaptive optimization-oriented redesign

Inputs: T, K, x_0 , prior $q(\Phi)$, state evolution and observation functions f, g

▷ For $t = 1$ to $T - 1$

Infer current state \hat{x}_t using u_{t-1}^*, \hat{x}_{t-1} and Thompson Sampling of prior $q(\Phi)$

▷ For $k = 1$ to K

$\hat{q}(\Phi_k) = q(\Phi)$

While $t < T$ (Forward Pass)

Thompson Sampling of the prior $\hat{q}(\Phi_k)$: $\phi_k = TS[\hat{q}(\Phi_k)]$

$u_t^k = \text{argmax}_{u \in \Omega} r_{t+1}(y_{t+1} | \phi_k, \hat{x}_t)$

Simulate redesign using u_t^k and predict $r_{t+1}(y_{t+1}^k | u_t^k, \hat{x}_t)$

Update prior: $\hat{q}(\Phi_k) \leftarrow \hat{q}(\Phi_k | u_t^k, y_{t+1}^k)$ using (u_t^k, y_{t+1}^k)

Accumulate reward: $\mathcal{R}_k = \mathcal{R}_k + r_{t+1}(y_{t+1}^k | u_t^k)$

End while

Define the policy: $\pi_t^k = \{u_t^k, \dots, u_{T-1}^k\}$ with its corresponding \mathcal{R}_k

▷ End for

Rank policies $\pi_t^k, k = 1, \dots, K$, using \mathcal{R}_k

Select the best policy $\pi_t^* = \{u_t^*, \dots, u_{T-1}^*\}$ with the highest \mathcal{R}_k

Redesign the experiment using only u_t^* and measure y_{t+1} at $t + 1$

Update prior: $q(\Phi) \leftarrow q(\Phi | u_t^*, y_{t+1})$ using experimental data (u_t^*, y_{t+1})

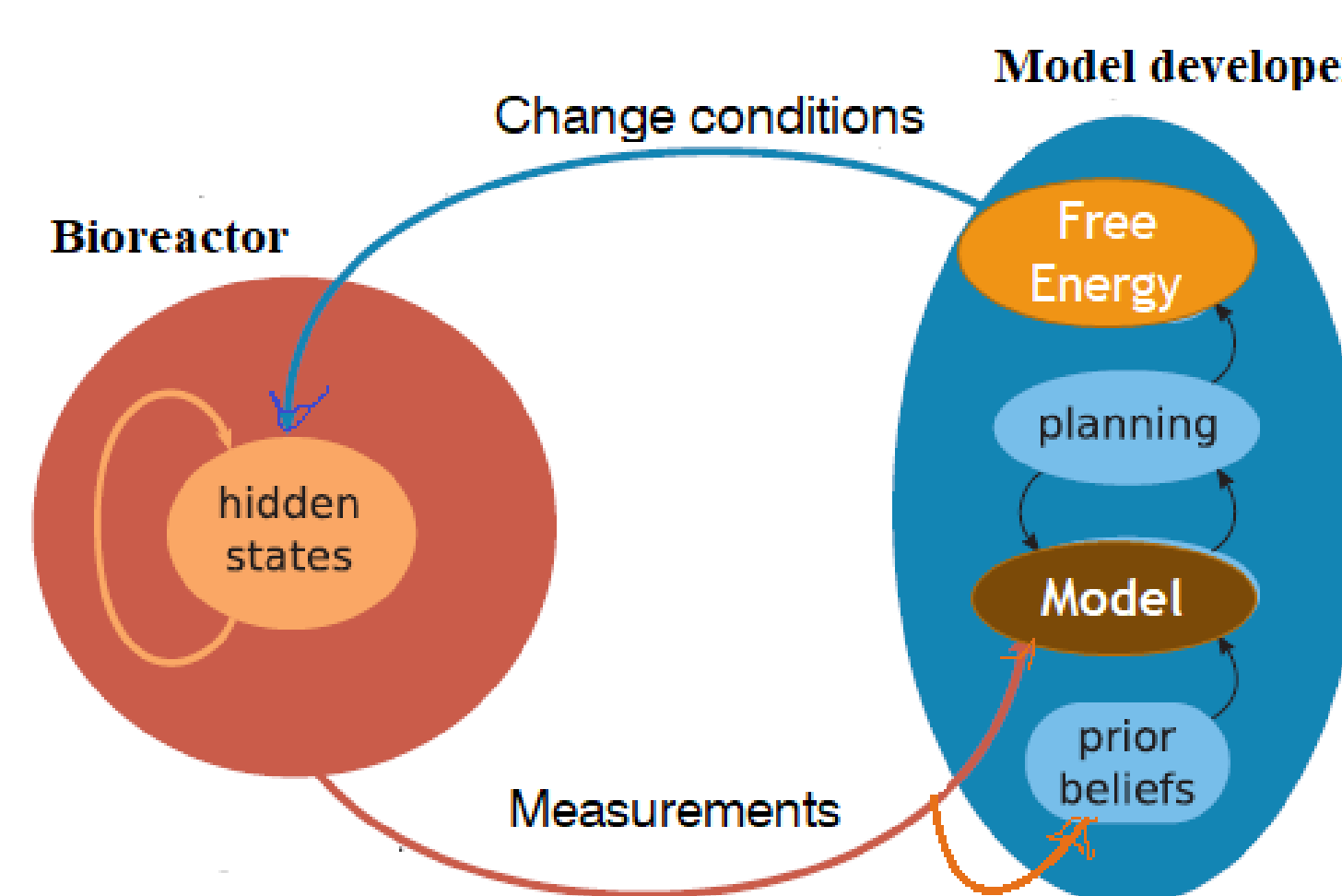
▷ End for

Outputs: $\pi^* = \{u_1^*, \dots, u_{T-1}^*\}, y = \{y_1, \dots, y_T\}, r = (r_1, \dots, r_T), q(\Phi)$

Active inference (Goal-directed sampling)

Active inference proposes that the modeler's goal or intent are encoded in the probabilistic model as a prior preference for desired observations (e. g., higher biomass productivity or protein expression).

Active Inference = Active Learning + Variational Inference



Probabilistic model learning is posed as the maximization of a free-energy lower bound $F(q)$ functional for the model evidence with respect to an approximate density q_ϕ

$$F(q) = \langle \ln p(\Phi | m) + p(\Phi | y, m) - p(\Phi) \rangle_q$$

$$ELBO = \ln p(y | m) - \mathcal{D}_{KL}(p(\Phi); p(\Phi | y, m))$$

where \mathcal{D}_{KL} is the **Kullback-Leibler divergence**.

Reinforcement learning for online redesign

Let $z_{t:T}$ denote a sequence of variables through time, $z_{t:T} = \{z_t, \dots, z_T\}$, and let define a **policy** Π as a **way of behaving** over time:

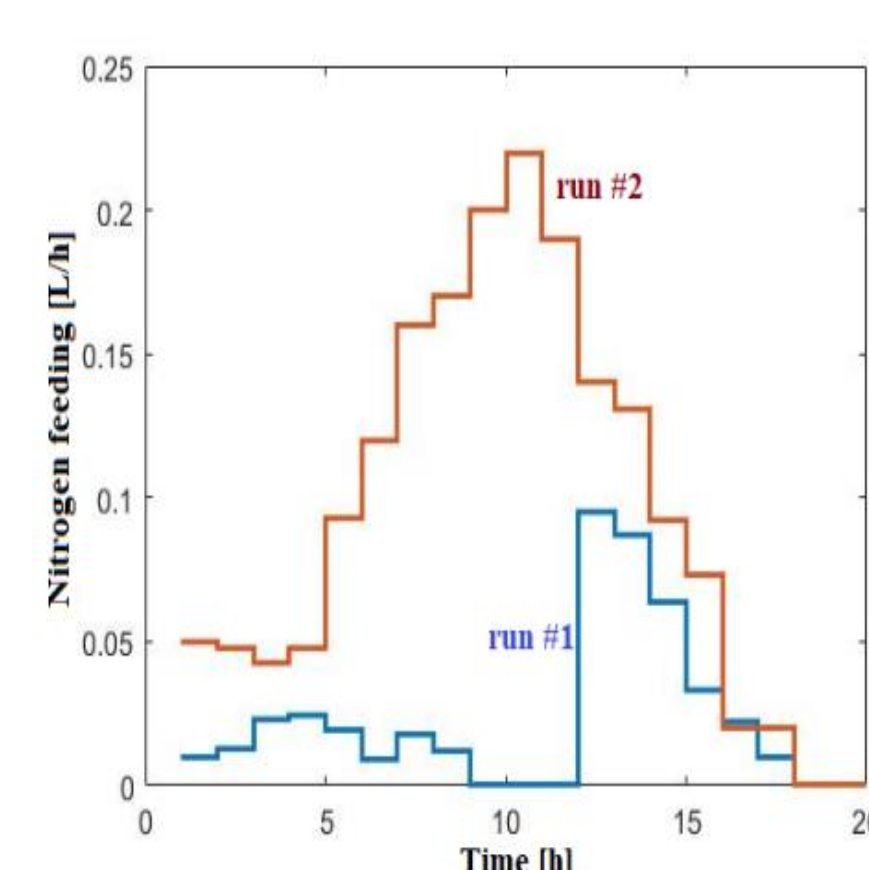
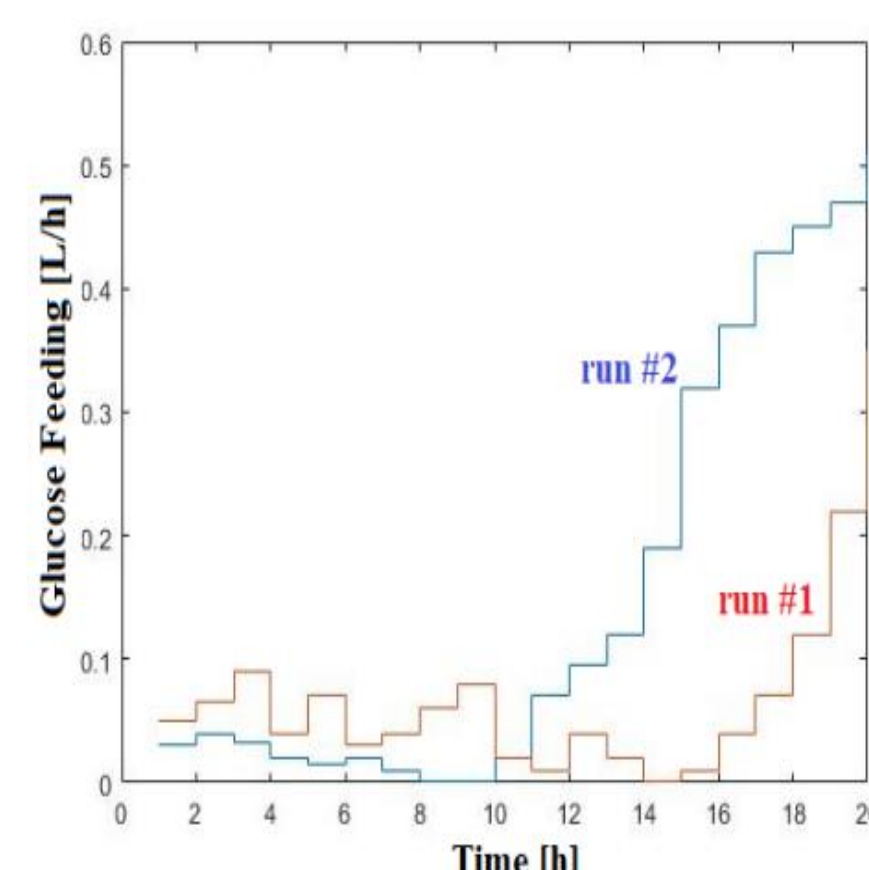
$u \leftarrow \Pi(x_t)$, (the map Π here is a probabilistic assignment from states to actions).

Applying recursively the policy Π defines **sequence of actions** $\pi = \{u_t, \dots, u_{T-1}\}$. In "modeling for optimization," the specific aim is to *minimize the free energy of the expected future* \tilde{F}_π , which is defined as:

$$\tilde{F}_\pi = \mathcal{D}_{KL}(q(y_{t:T}, x_{t:T}, \Phi | \pi) \| p^*(y_{t:T}, x_{t:T}, \Phi)); \quad \Phi = (\theta, \varphi, \alpha, \sigma)$$

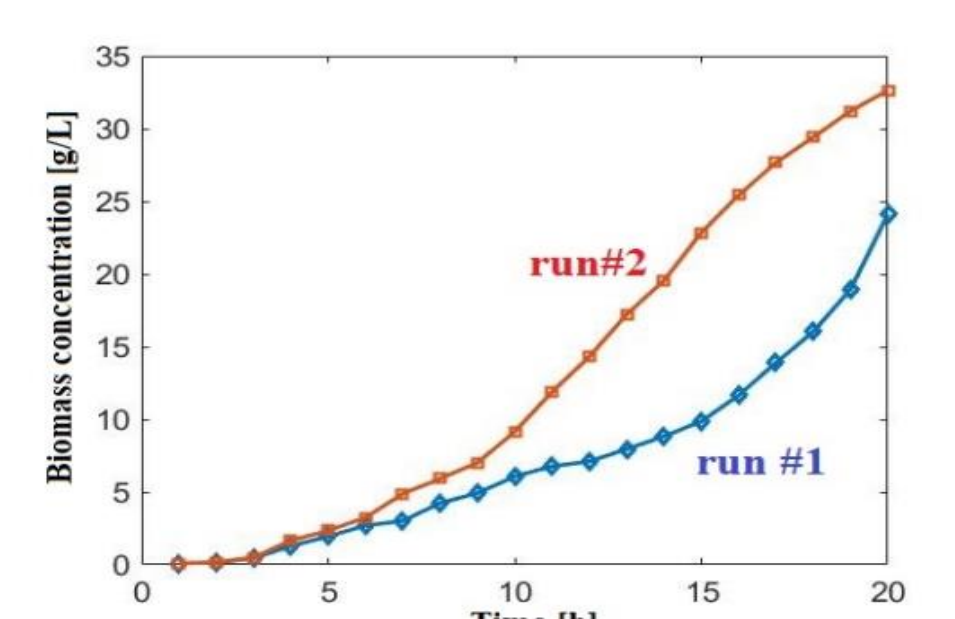
where $q(y_{t:T}, x_{t:T}, \Phi | \pi)$ models the probability distribution for future trajectories in a dynamic experiment under policy π and $p^*(y_{t:T}, x_{t:T}, \Phi)$ defines the joint probability distribution for the optimal trajectory of the hidden states, model parameters and preferred observations. Thus, when \tilde{F}_π is driven to zero, the policy π becomes the (probabilistic) optimal policy. Notice that by minimizing \tilde{F}_π , the surprise $-\ln p(y_{t:T} | m)$ is also minimized, which in turns maximize the Bayesian model evidence.

Results for the Baker's yeast production process



Units	Run #1	Run #2	Run #3
μ	0.5431	0.250	0.6232
σ	0.0612	0.020	0.0578
μ	0.8929	0.085	0.8500
σ	0.2647	0.080	0.2450
μ	0.2589	0.080	0.0189
σ	1.0150	0.150	0.9733
μ	0.4445	0.125	0.4210
σ	2.5364	0.200	2.6472
μ	1.1903	0.150	1.2279
σ	0.1524	0.030	0.1208
μ	3.1817	0.050	3.2011
σ	2.9370	0.050	2.9674
μ	9.0014	2.000	9.4569
σ	5.9981	0.500	5.8919
μ	5.7377	0.210	6.1311

Modeling Run #	Biomass [g/L]
1	24.11
2	32.64
3	31.58
Richelle et al. 2014	32.00



Significance and Concluding Remarks

A novel probabilistic method for modeling the dynamic behavior of bioreactors in the most profitable region of operating conditions is proposed. Based on simulation data, a dynamic experiment is redesigned online through active inference. Reinforcement learning is used to maximize the Bayesian model evidence, that is, to minimize surprise.

- ✓ Probabilistic (causal) models of bioreactors are learned by biasing data gathering using the **Free Energy of the Expected Future**.
- ✓ Reinforcement learning following an MPC-approach is used to **combine planning and control for online experiment redesign**.
- ✓ Bayesian Variational Analysis methods are applied for state inference and probabilistic parameter estimation.
- ✓ Simulation data is used to **learn a redesign policy** for adaptive experimental design.

Acknowledgements / References