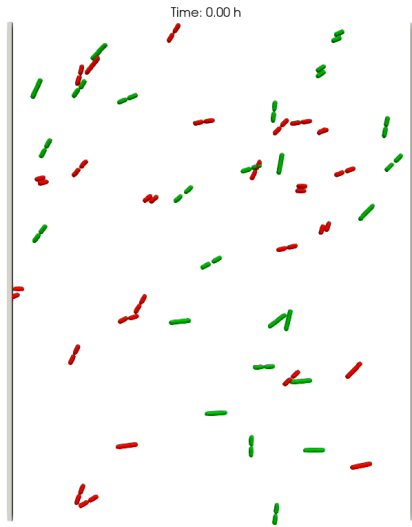
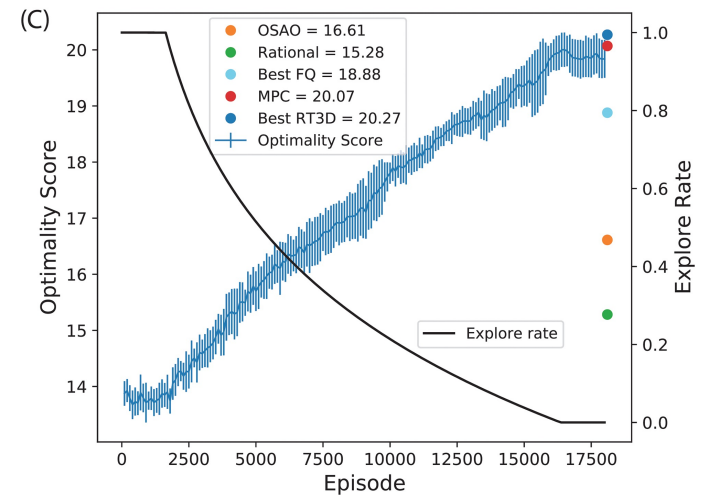
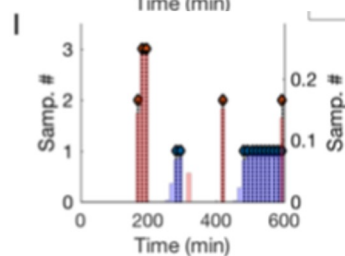
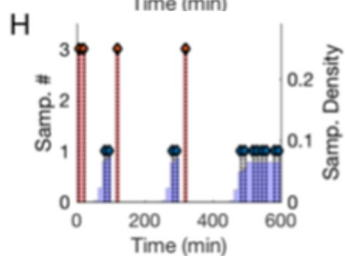
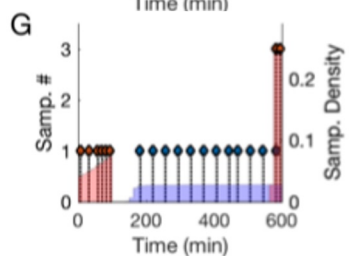
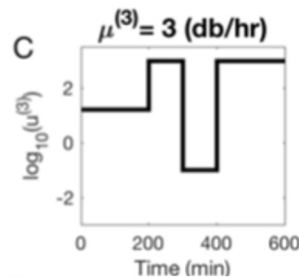
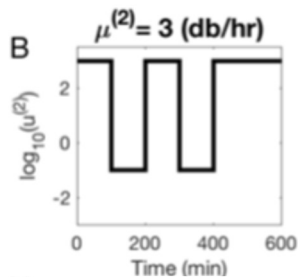
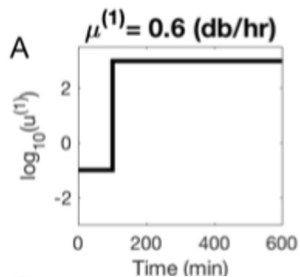
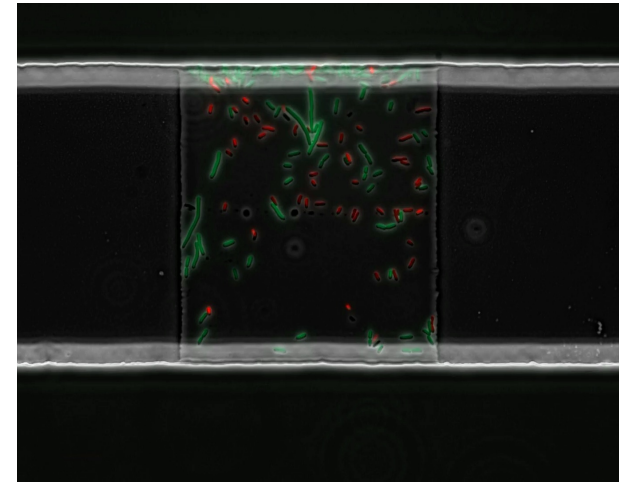


Challenges and opportunities for model calibration



Brian Ingalls (he/him)
 Department of
 Applied Mathematics
 University of Waterloo
 Waterloo, Canada
 @bpingalls



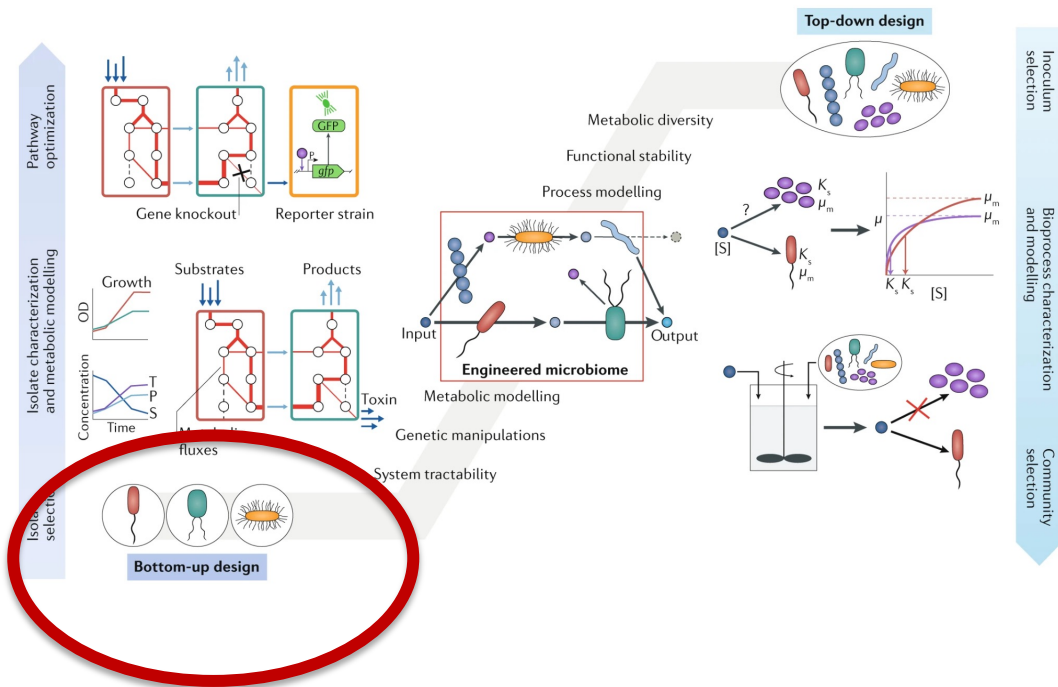
Outline

- 1) Calibration strategies for agent-based population models of mixed bacterial populations
- 2) Optimal experimental design tools for systems and synthetic biology

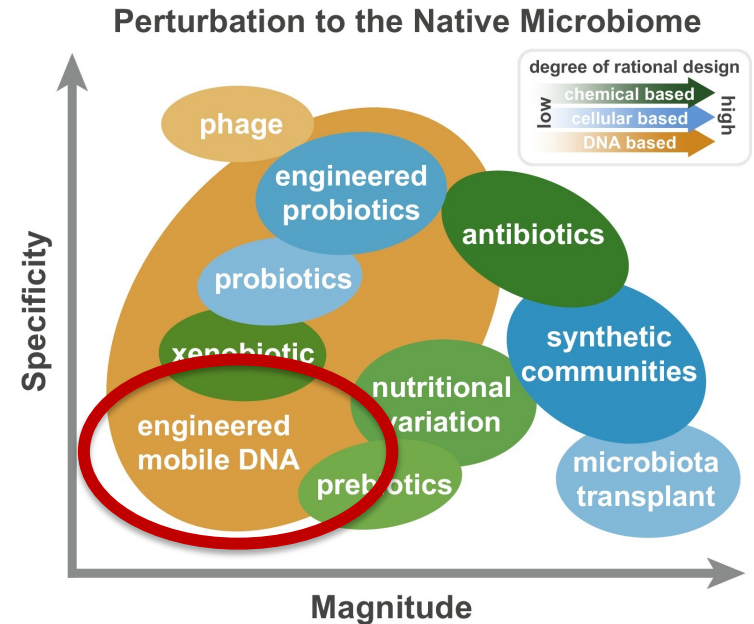
Outline

- 1) Calibration strategies for agent-based population models of mixed bacterial populations**
- 2) Optimal experimental design tools for systems and synthetic biology

Goal: model-based design for manipulation of mixed microbial communities

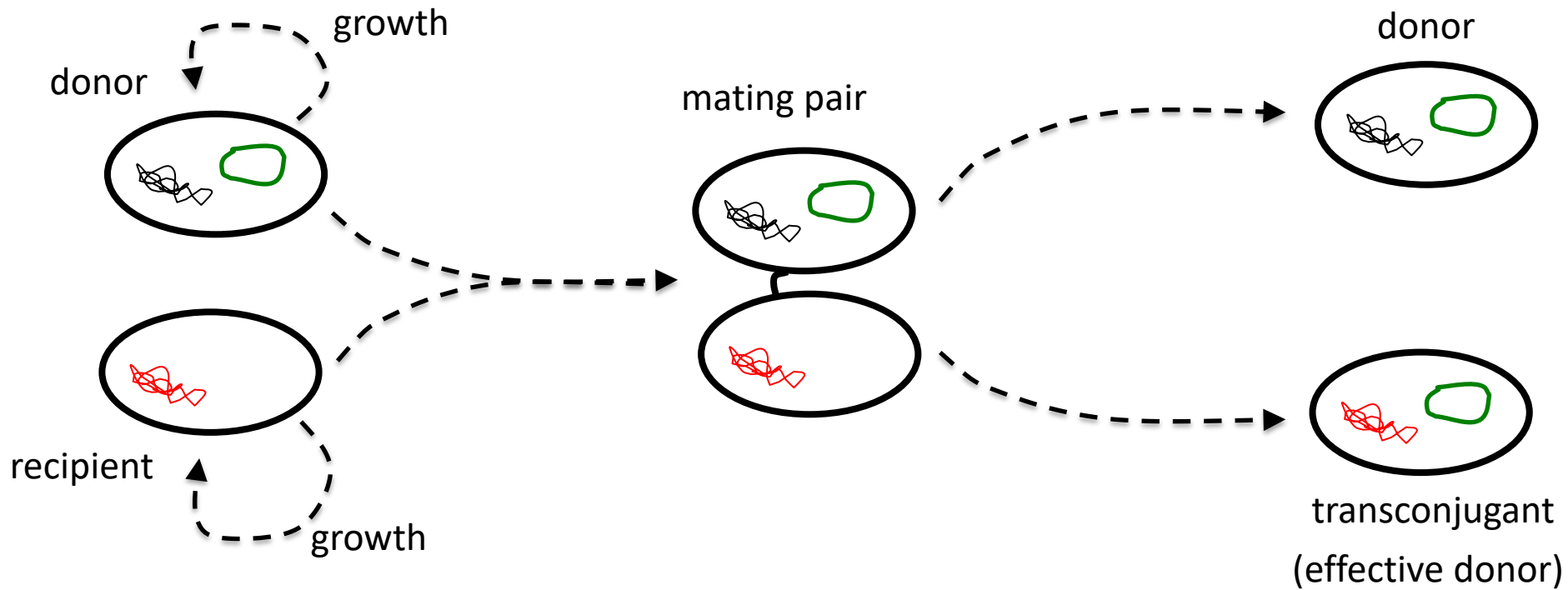


Lawson, et al, *Nature Reviews Microbiology*, 2019

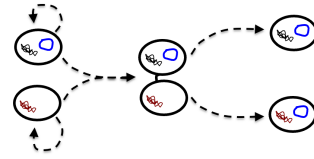
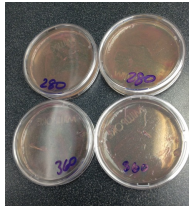


Sheth, et al. *Trends in Genetics*, 2016

Modelling of plasmid delivery by conjugation



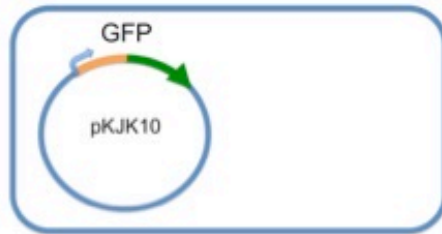
Approach 1: Filter Mating Experiments



Populations binned by fluorescence signature

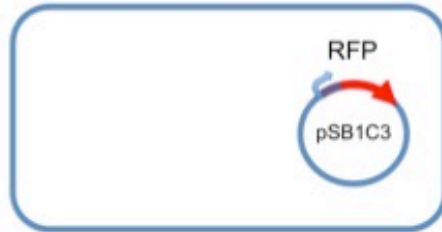
E. coli CSH26

Donors



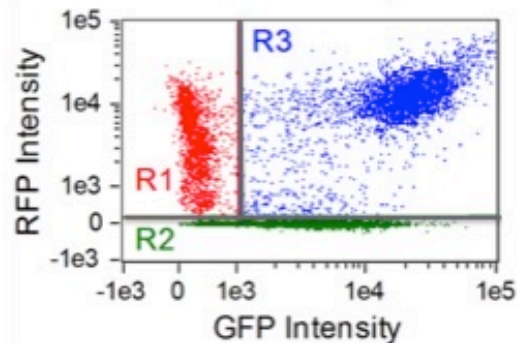
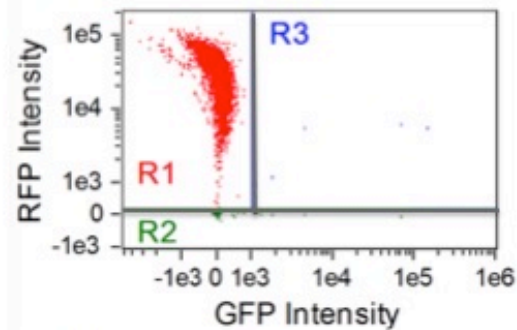
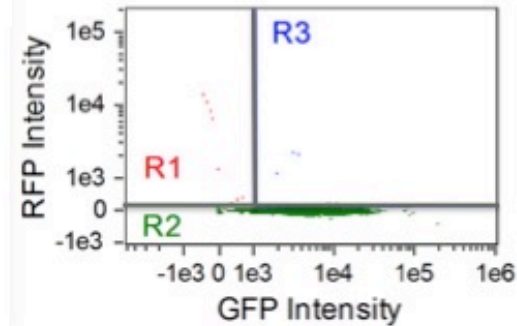
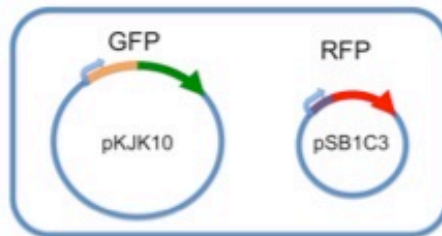
E. coli DH5 α

Recipients

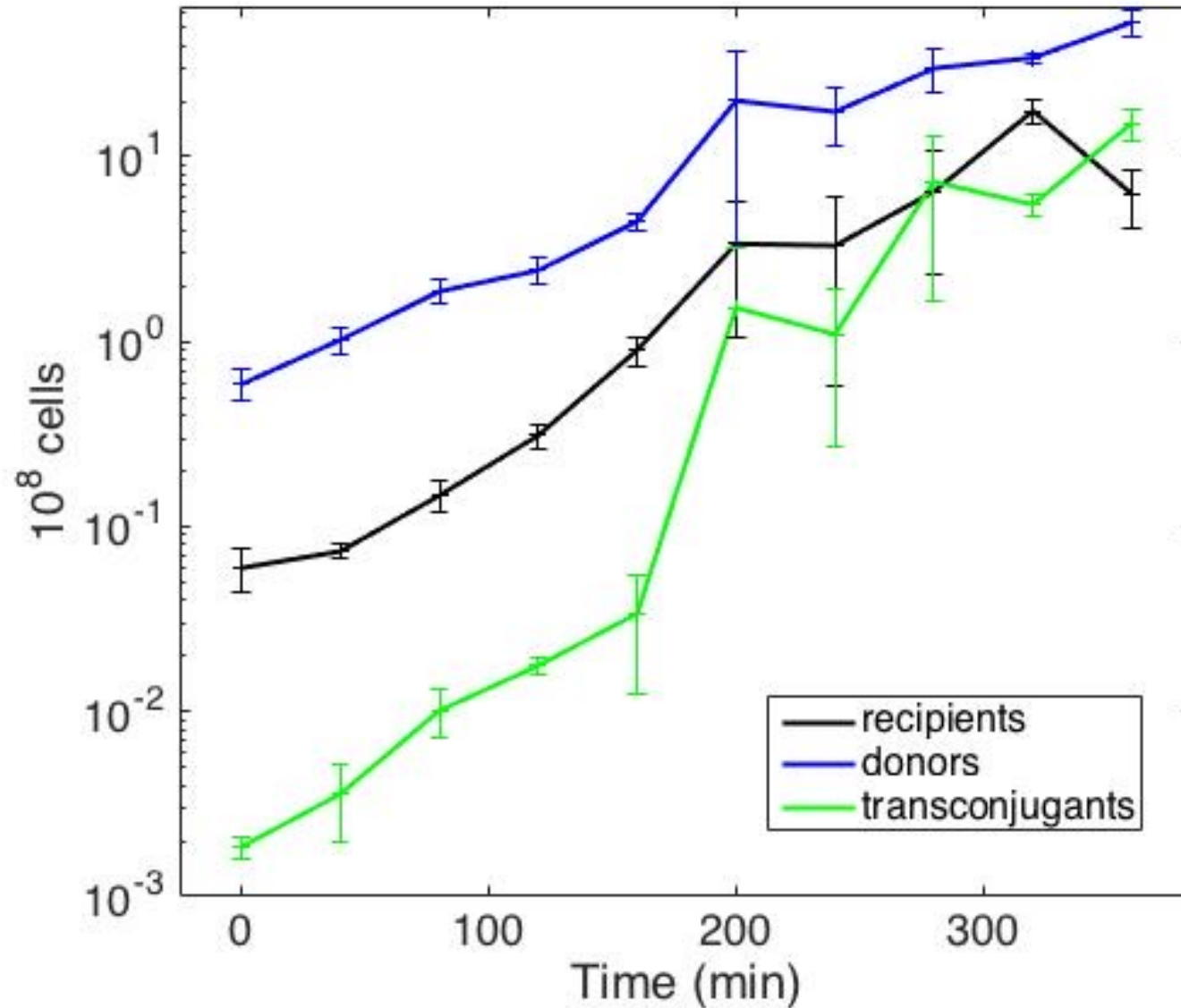


E. coli DH5 α

Transconjugants



Time point observations



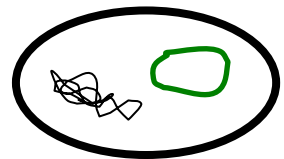
Differential Equation Model (Levin et al., 1979)

donor population: D

recipient population: R

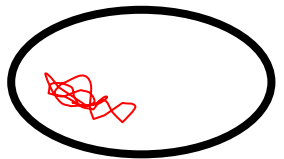
transconjugant population: T

Balance equations:



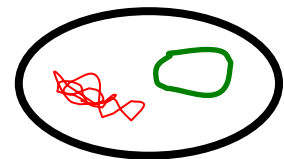
Donors

$$\frac{d}{dt}D(t) = \alpha D(t)_{\text{growth}}$$



Recipients

$$\frac{d}{dt}R(t) = \alpha R(t)_{\text{growth}} - \underbrace{\gamma(D(t) + T(t))}_{\text{effective donors}} R(t)_{\text{conjugation}}$$

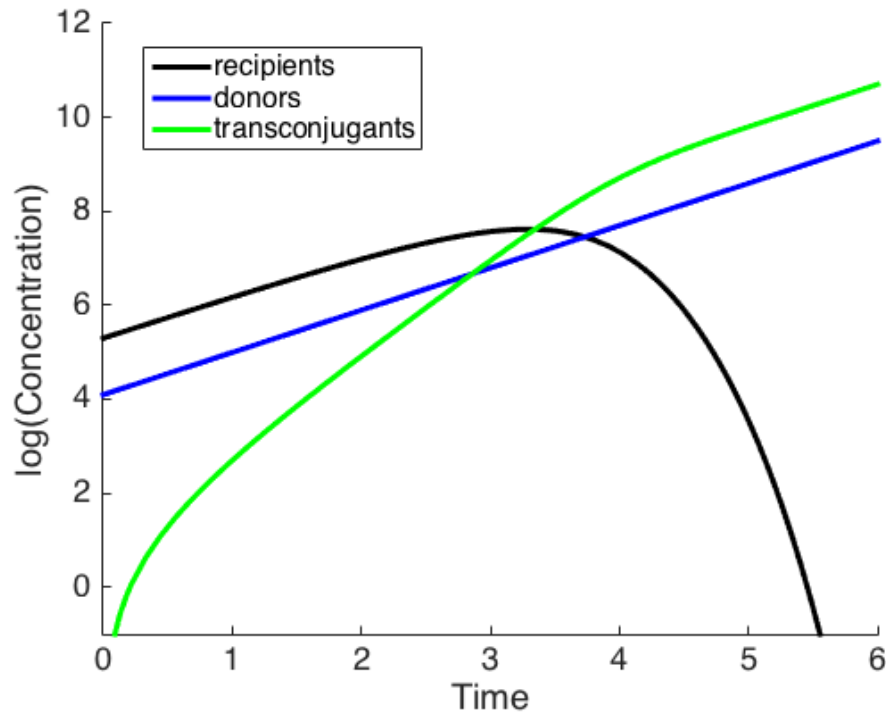


Transconjugants

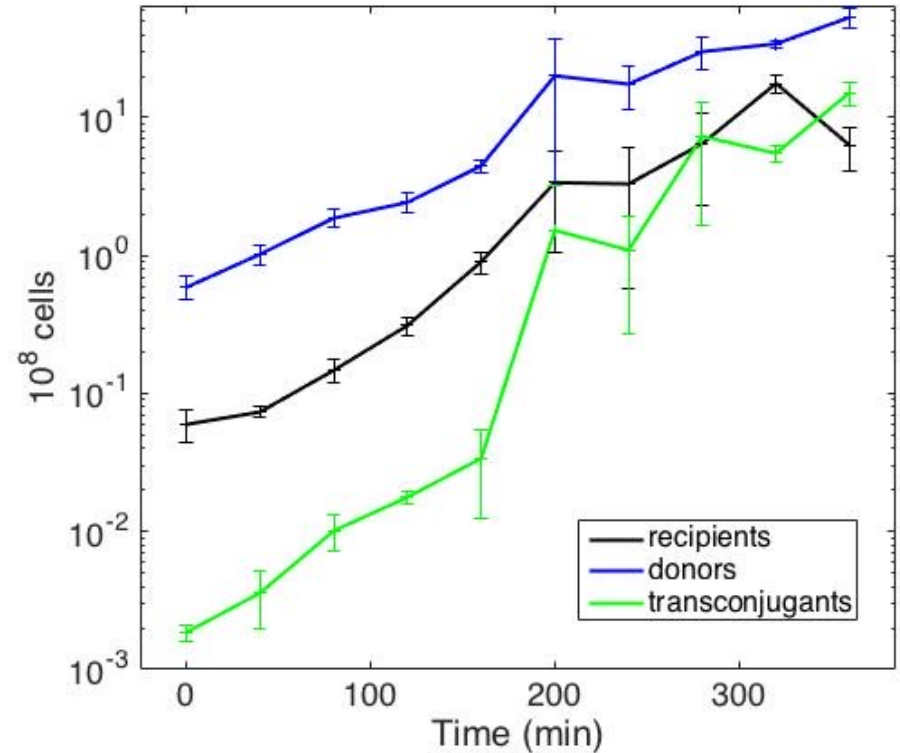
$$\frac{d}{dt}T(t) = \alpha T(t)_{\text{growth}} + \underbrace{\gamma(D(t) + T(t))}_{\text{conjugation}} R(t)$$

Analogous to susceptible-infectious (SI) epidemiological models

Levin model



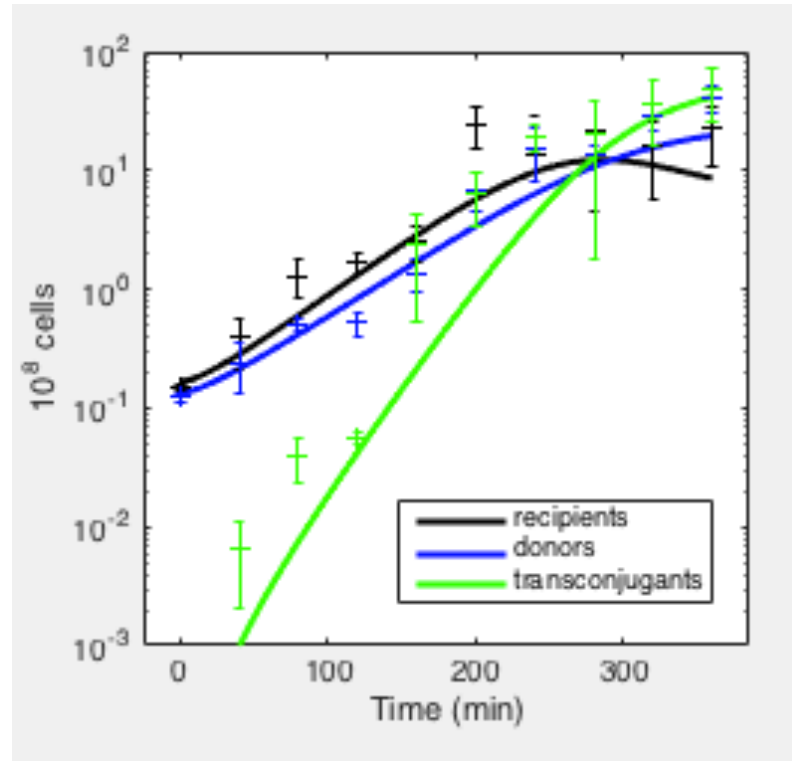
Filter mating data



Adjustments required:

- Distinct kinetics for donors, recipients, and transconjugants
- Lag in initial growth (lag phase) [Baranyi and Roberts, 1994]
- Nutrient limitation (stationary phase) [Simonsen et al., 1990]

Parameter Fitting

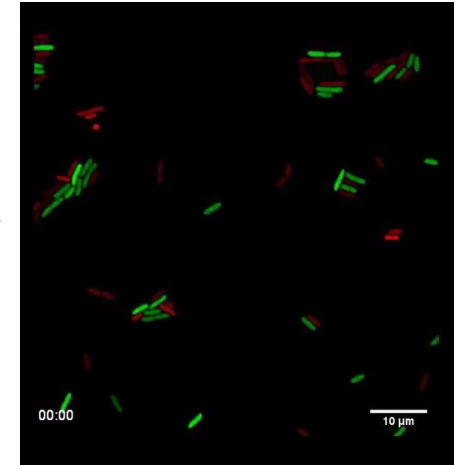
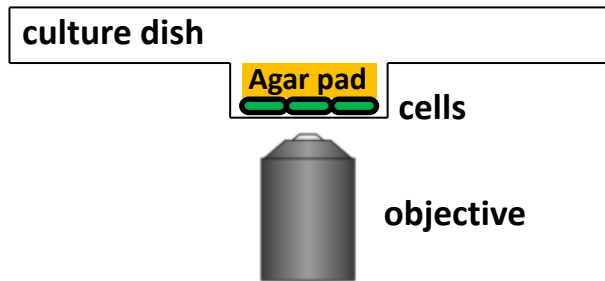


Quality of fit: weighted sum of squared errors

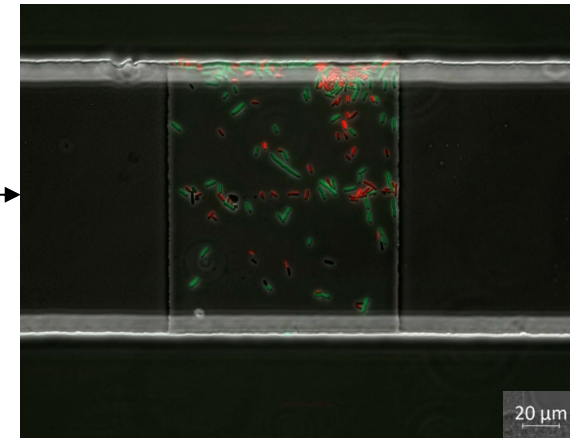
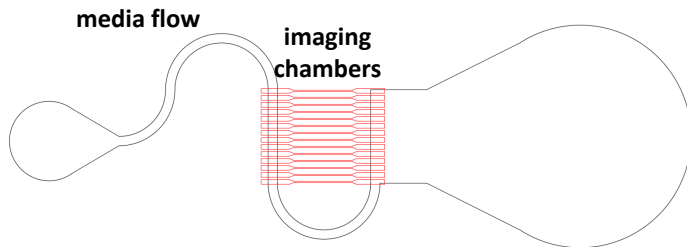
$$\text{Error}(\mathbf{p}) = \sum \left(\frac{\text{observation} - \text{prediction}}{\text{standard deviation}} \right)^2$$

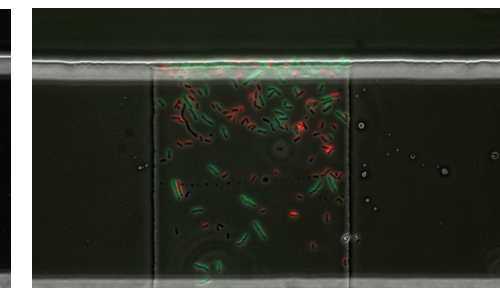
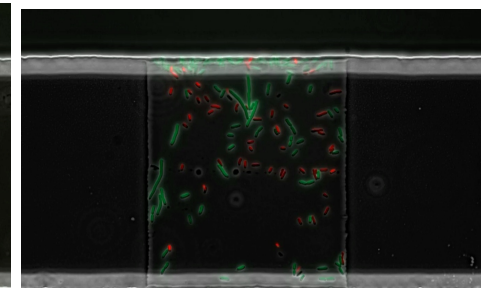
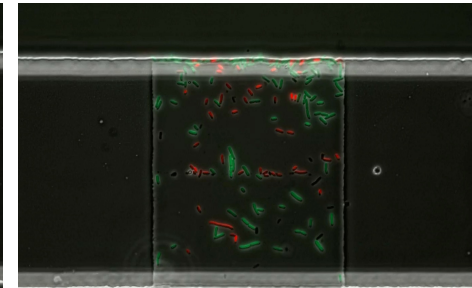
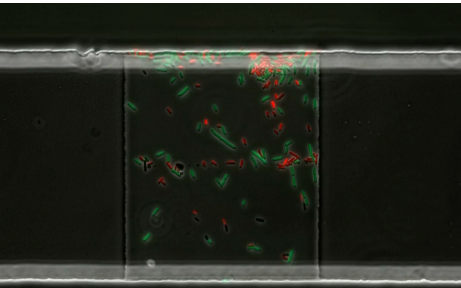
Approach 2: Collection of spatiotemporal data

Agar pad



Microfluidic chip





Microfluidics

experiment: 38 h

frequency: 6 mins

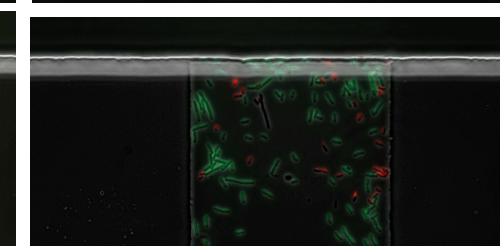
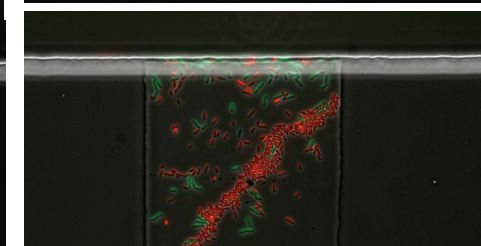
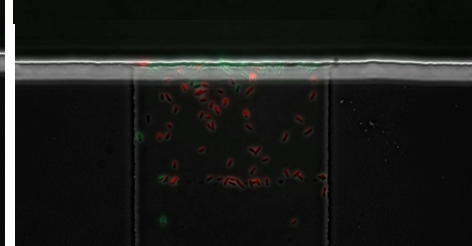
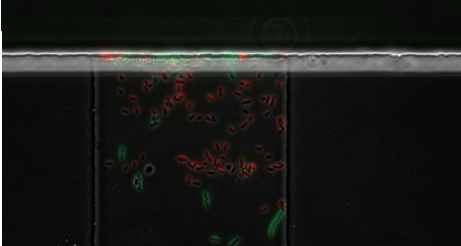
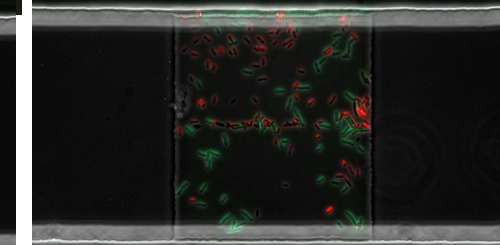
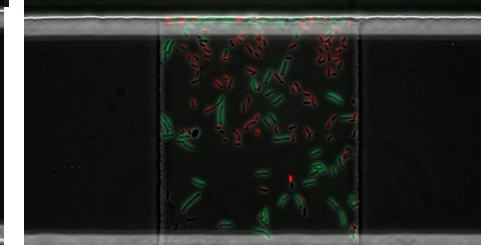
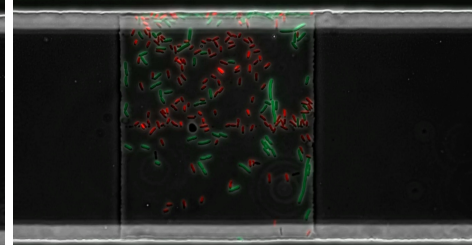
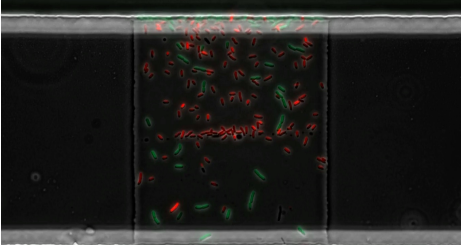
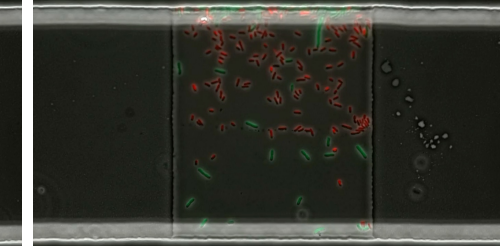
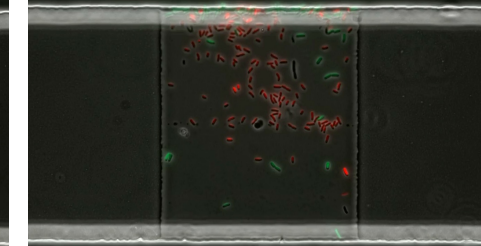
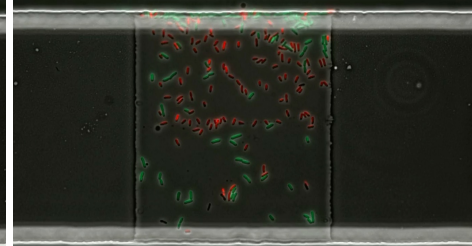
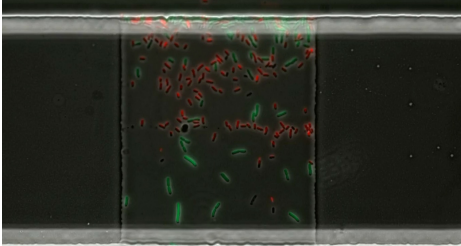
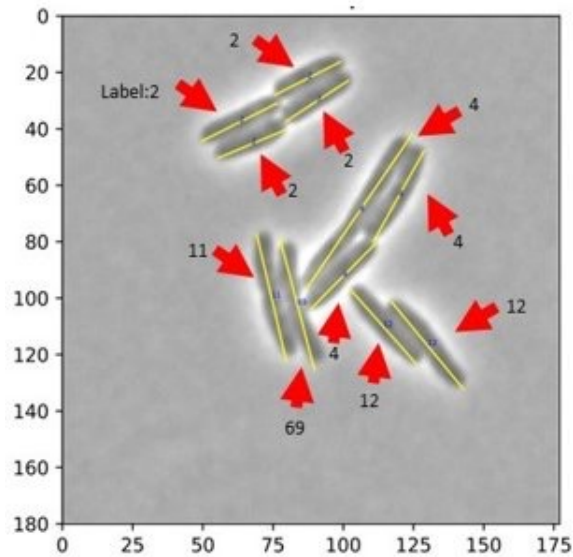
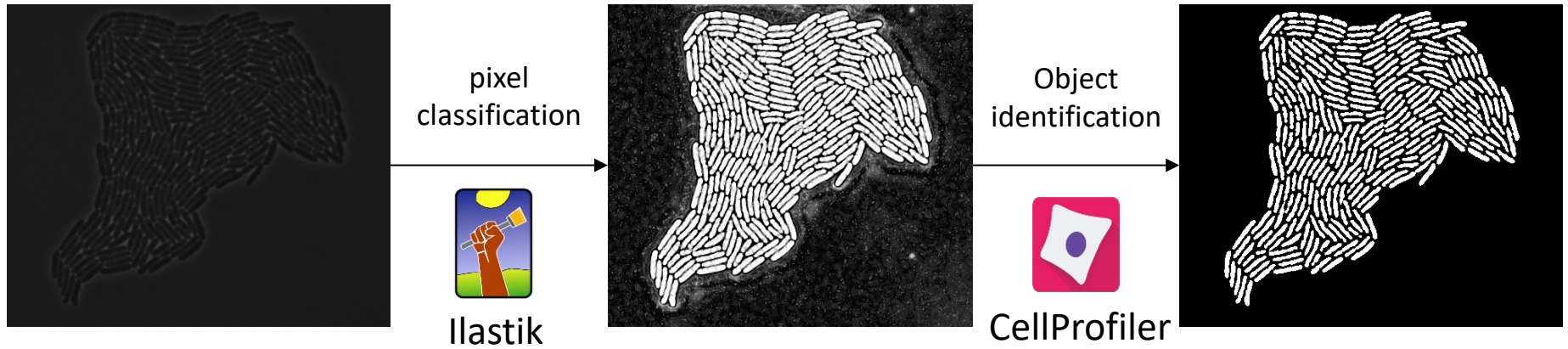
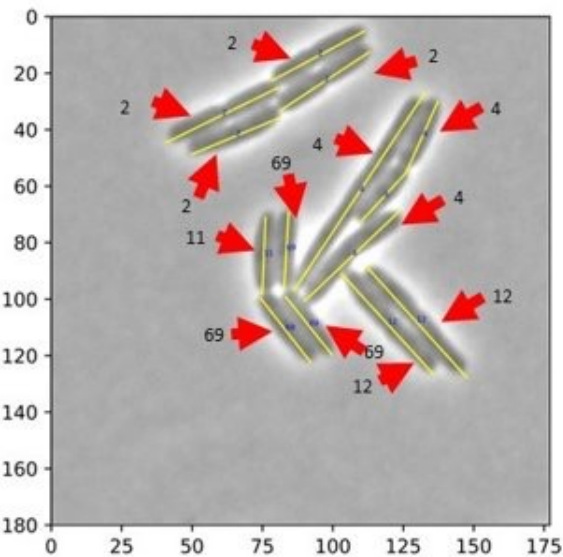


Image processing

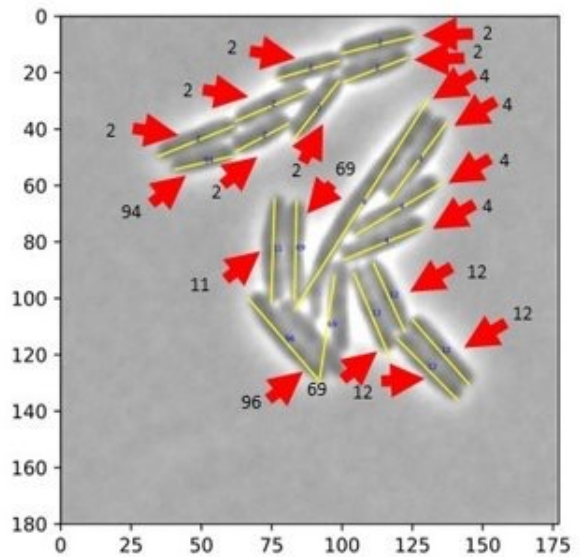
Segmentation



$t=45(\text{min})$



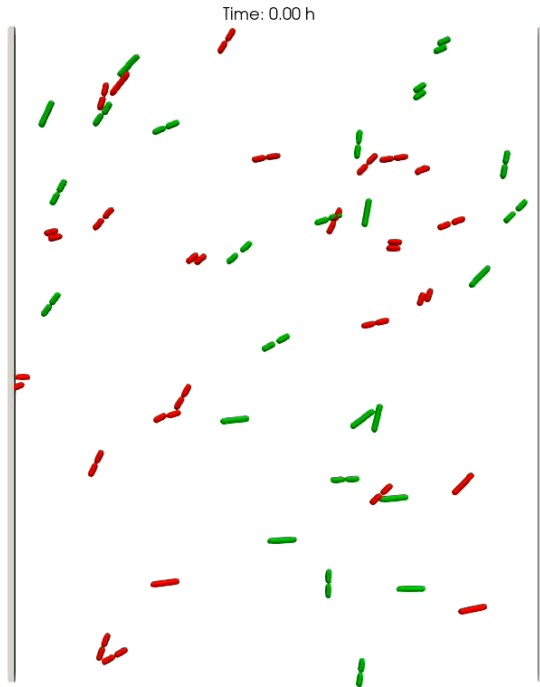
$t=75(\text{min})$



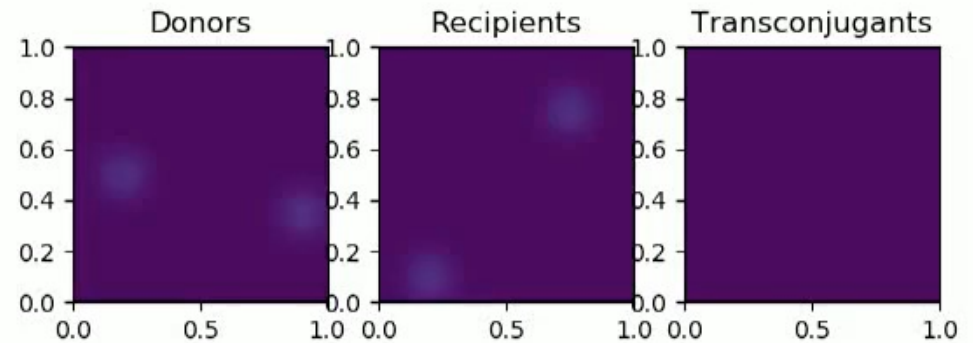
$t=105(\text{min})$

Frame-to-frame: track cells and identify division events

Modelling approaches



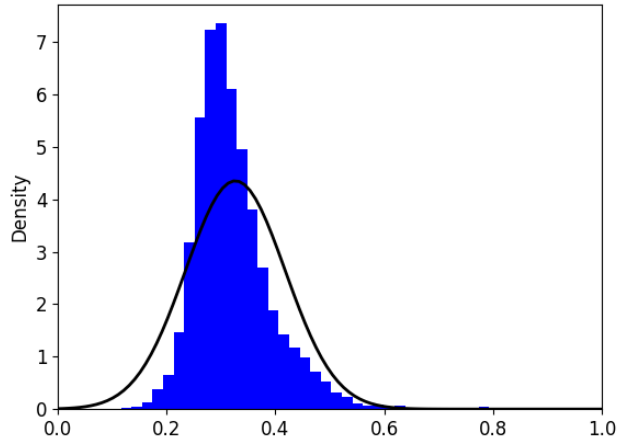
Single-cell: **Individual/Agent-based model**



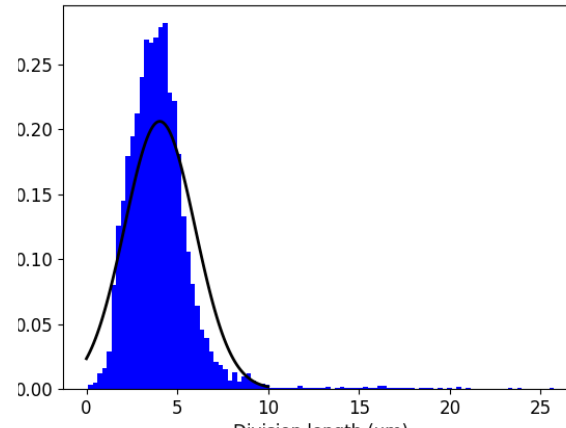
Coarse-grained (density measure):
partial differential equation

Directly observable parameters

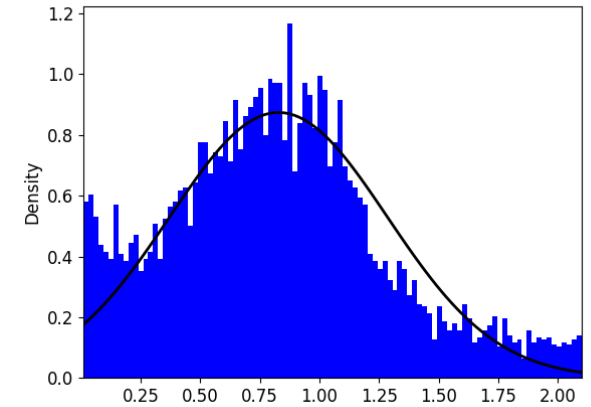
Individual cellular measurements



Radius (um)

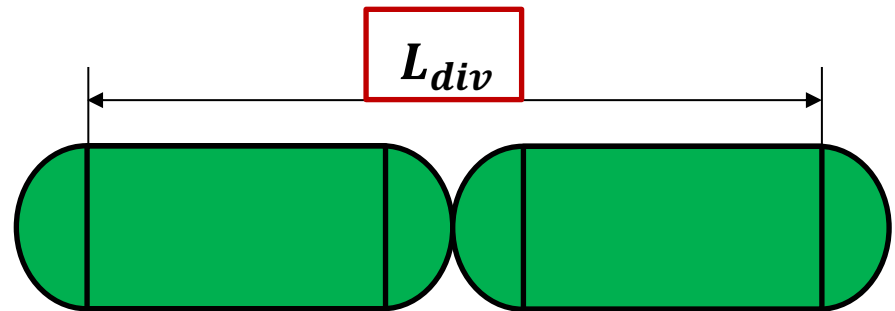
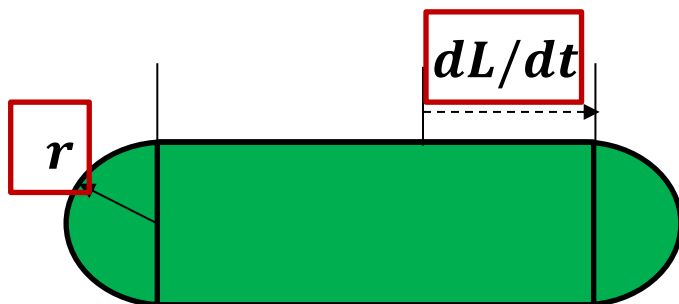


Length at division (um)



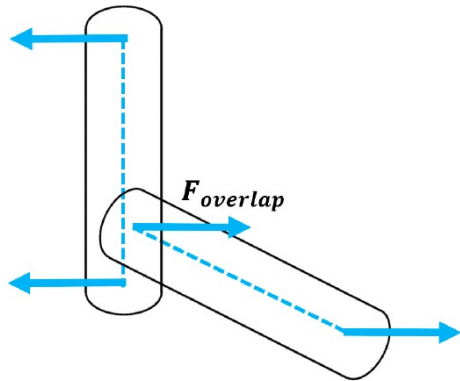
Elongation rate (um/h)

These can be incorporated directly into the ABM formulation

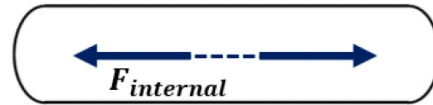


Parameters to be inferred

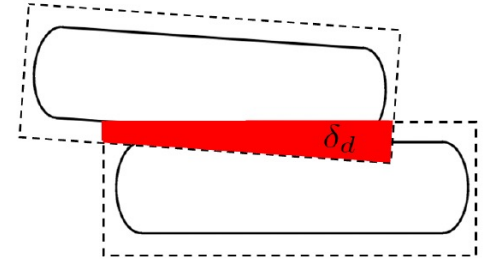
Biophysical:



cell stiffness
(exclusion force)

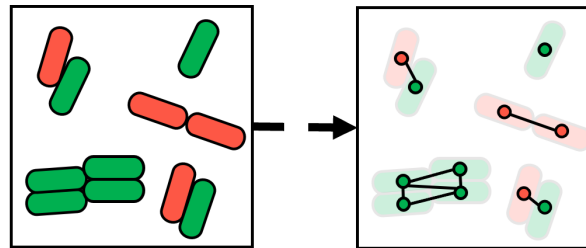


growth pressure



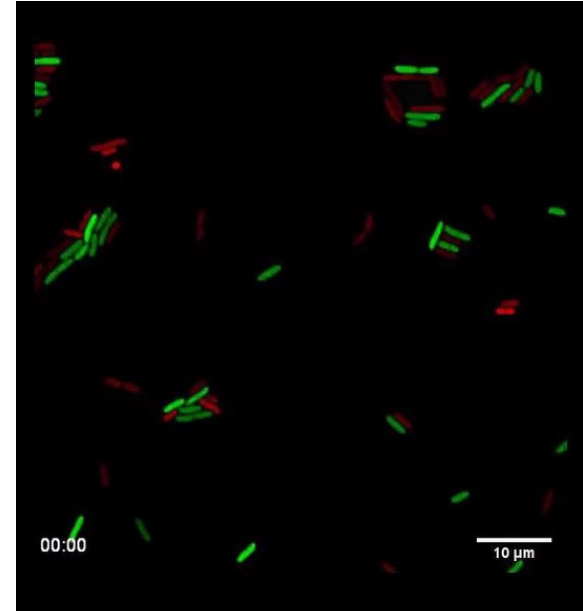
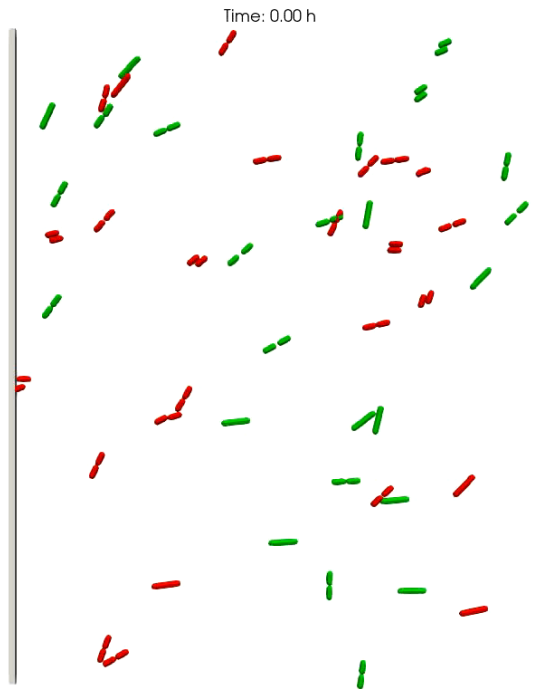
adhesion strength

Process-specific:



conjugation process
(degree of contact, delay, zygotic induction)

Agent-based model calibration challenges

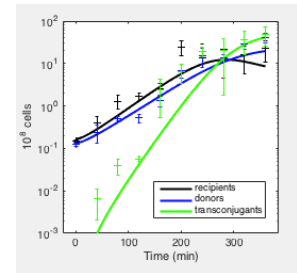


Stochasticity

Lack of 'obvious' goodness-of-fit measure

Strategy (from ecology):

“Pattern-oriented Modelling”

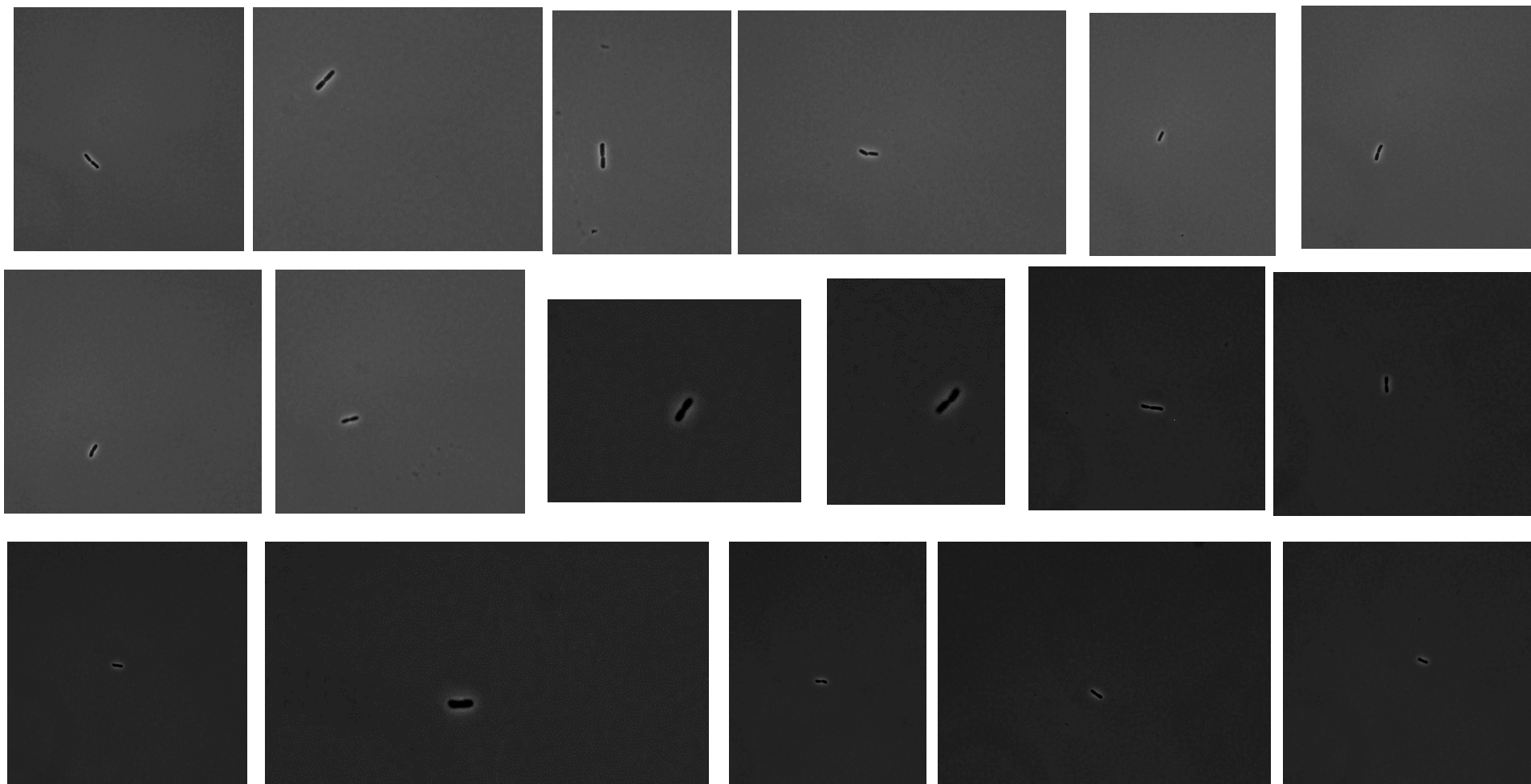
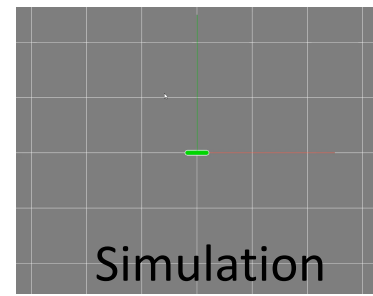


$$\text{Error}(\mathbf{p}) = \sum \left(\frac{\text{observation} - \text{prediction}}{\text{standard deviation}} \right)^2$$

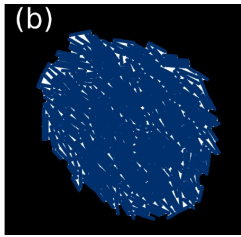
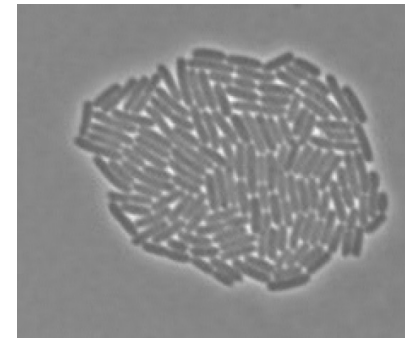
non-spatial case

Observation of system behavior

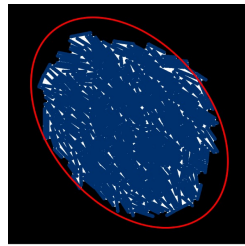
Biophysical parameters: growth of isolated microcolonies



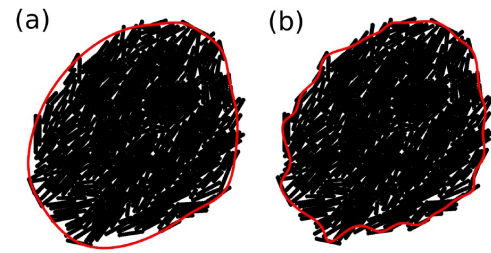
Biophysical Features



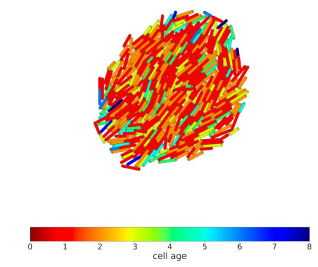
colony density



colony aspect ratio



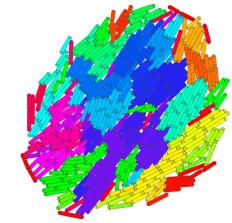
boundary jaggedness



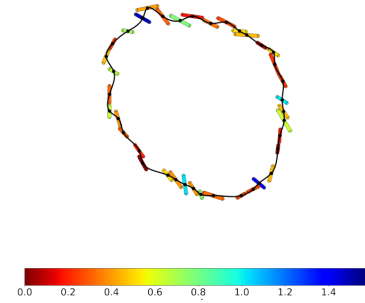
distribution of cell (pole) ages



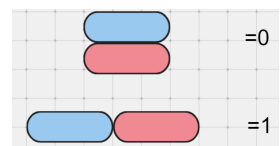
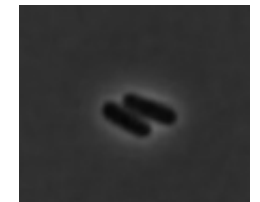
anisotropy of cell orientation



distribution of oriented 'patches'



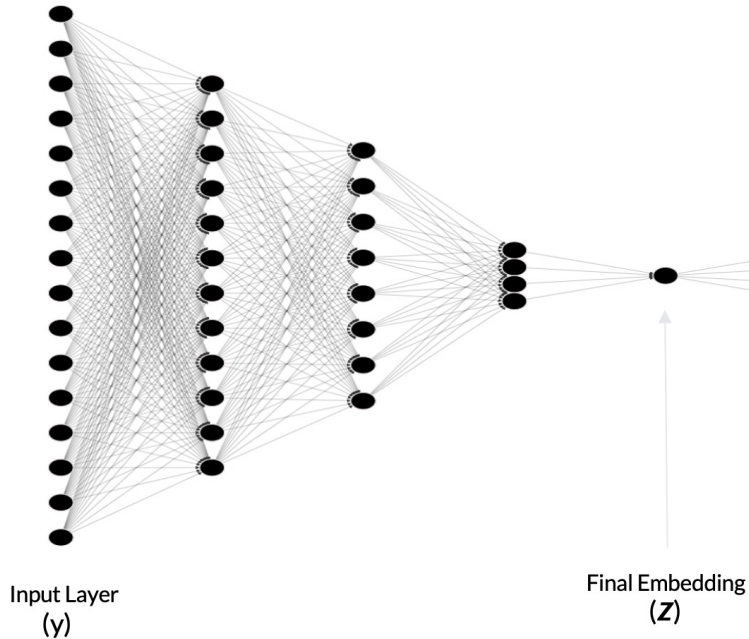
orientation of boundary cells



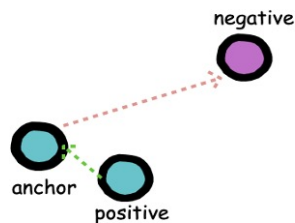
initial dyad orientation

Deep representation learning to identify features

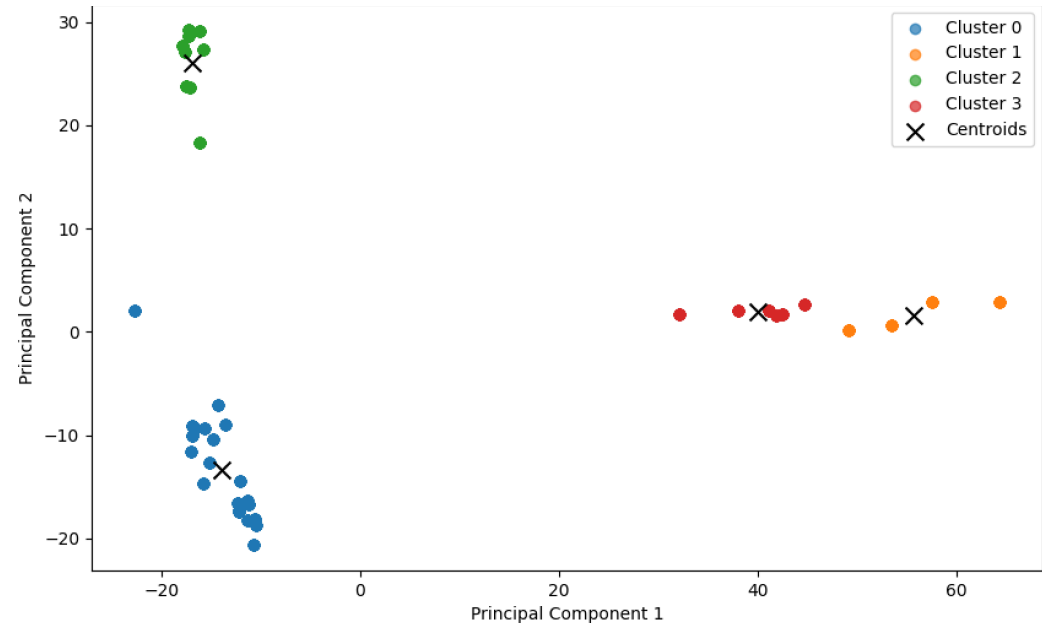
Inspired by: Cess and Finley "Calibrating agent-based models to tumor images using representation learning." *PLOS Comp. Biol.* 2023



Layers: LSTM or convolutional
Triplet loss (supervised learning)

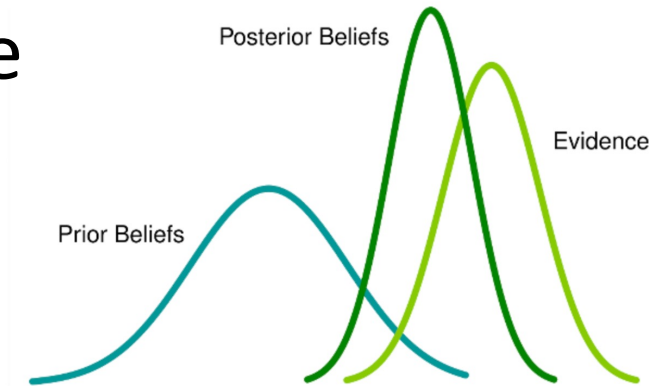


Successful preliminary results with low-dimensional embeddings (simulated data)



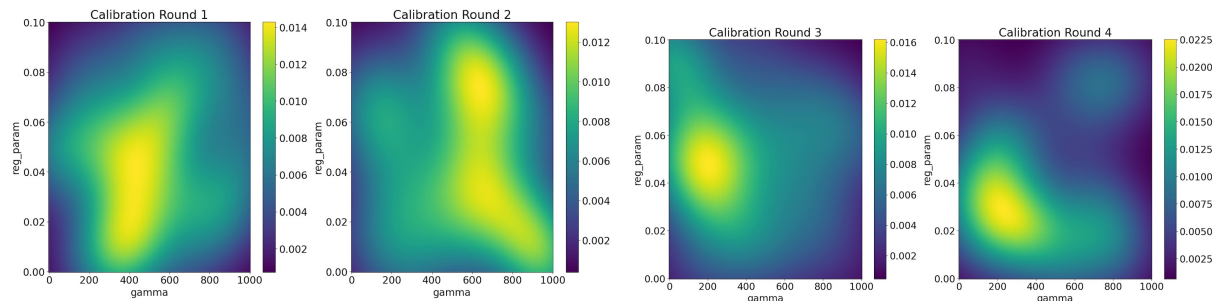
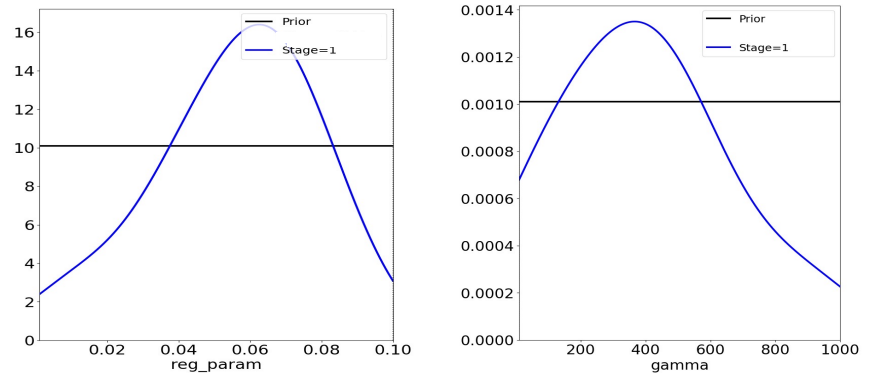
Calibration procedure

Approach: Bayesian inference



Method: Approximate
Bayesian Computation
(ABC)

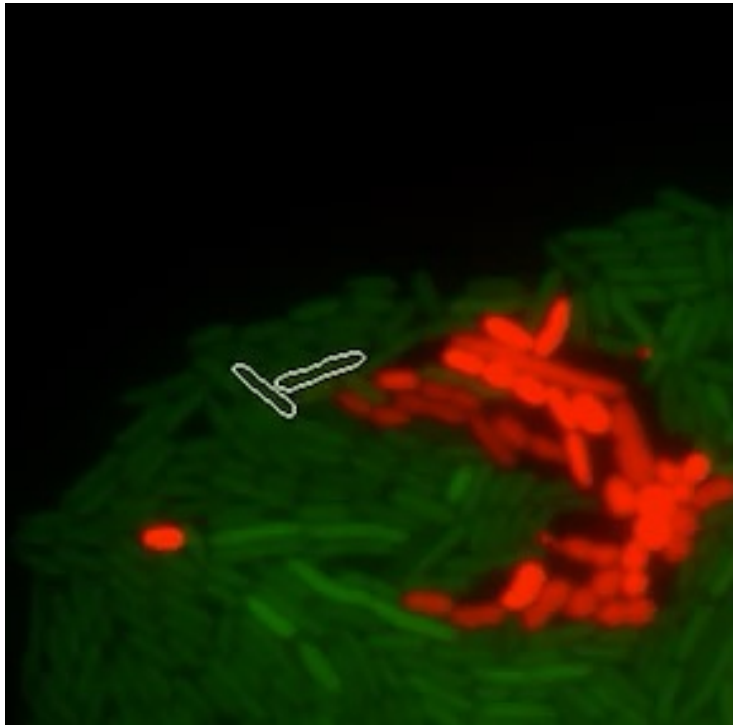
Results:



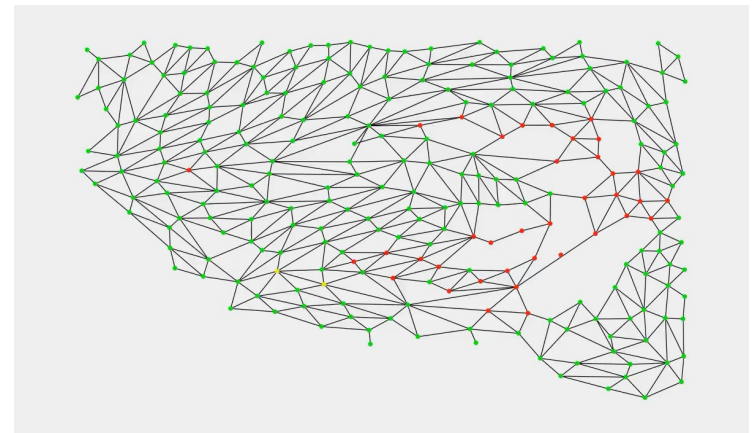
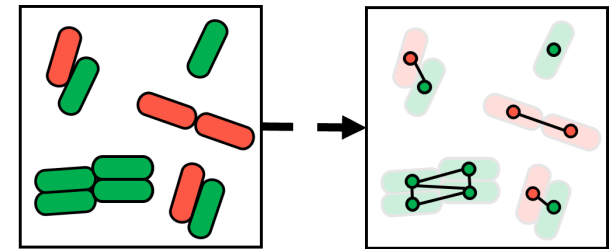
Also exploring deep regression models

Next step: calibration of conjugation parameters (incubation period*, degree of contact, zygotic induction)

Delayed conjugation events



Contact network



Identification of conjugation events: integer programming approach inspired by epidemiological contact tracing analysis

Outline

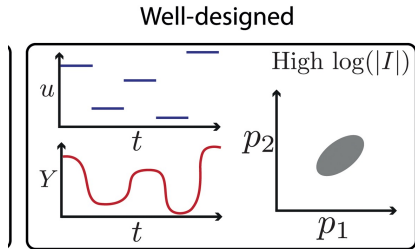
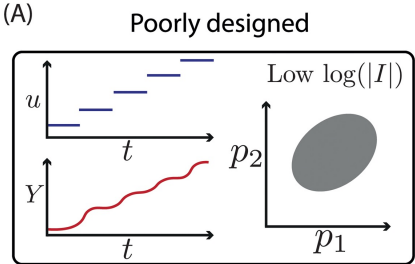
- 1) Calibration strategies for agent-based population models of mixed bacterial populations
- 2) **Optimal experimental design tools for systems and synthetic biology**

Model Based Optimal Experimental Design

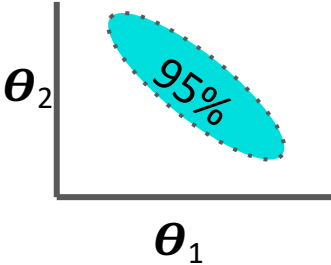
Identification of maximally informative experiments

Fisher Information-based approach

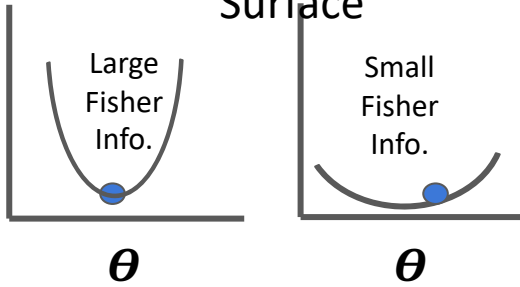
D-optimality



Minimize the Confidence Region Volume



Find an Experiment that Yields a 'Pointy' Optimization Surface



$$\text{cov}(\theta) \sim I^{-1}$$

$$\det(\text{cov}(\theta))$$

Minimize **D**eterminant of the Parameter Covariance Matrix

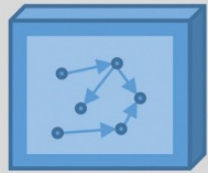
$$\det(I)$$

Maximize **D**eterminant of Fisher Information Matrix*

*Sensitivities of outputs to parameter values scaled by confidence in measurements

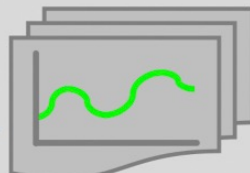
Tools for model-based dynamic experimental design

Candidate models

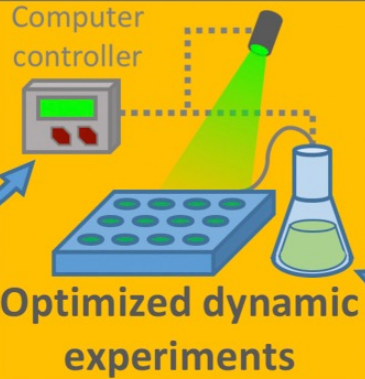


Dynamic models

MBDOE algorithm



Input signal

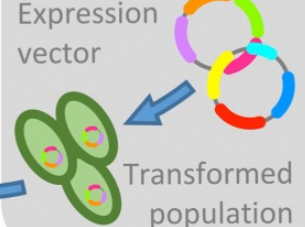


Tools for dynamic biological experiments

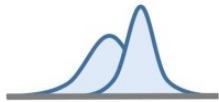
Experimental apparatus



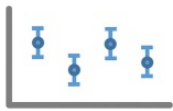
Dynamic control system



Estimate parameters



Improve predictions



Select model structure

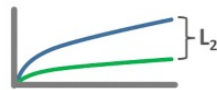


Modelling goals

Fisher information matrix methods

$$\begin{bmatrix} \frac{\partial y_1}{\partial \theta_1} & \frac{\partial y_2}{\partial \theta_1} \\ \frac{\partial y_1}{\partial \theta_2} & \frac{\partial y_2}{\partial \theta_2} \end{bmatrix}$$

Model discrimination



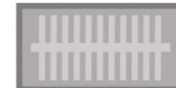
Bayesian methods

$$P(\theta|x) \propto P(x|\theta)P(\theta)$$

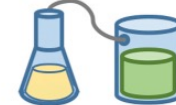


Optimization criteria

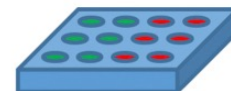
Microfluidics



Continuous culture



Optical array



Experimental hardware

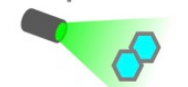
Chemically induced gene expression



Optically induced gene expression



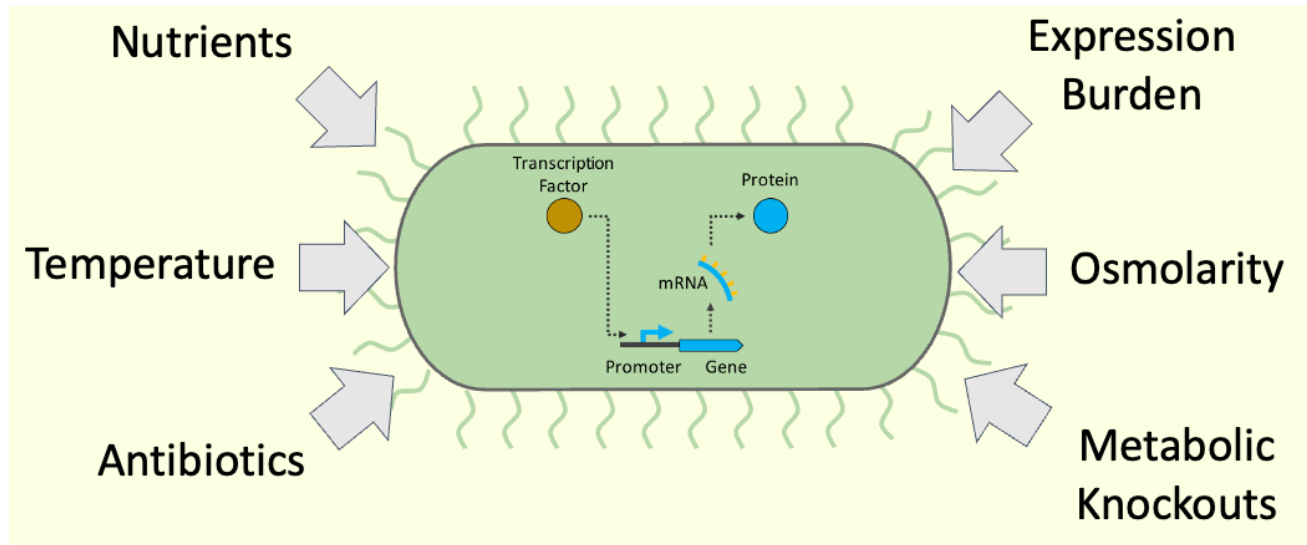
Photo-activated proteins



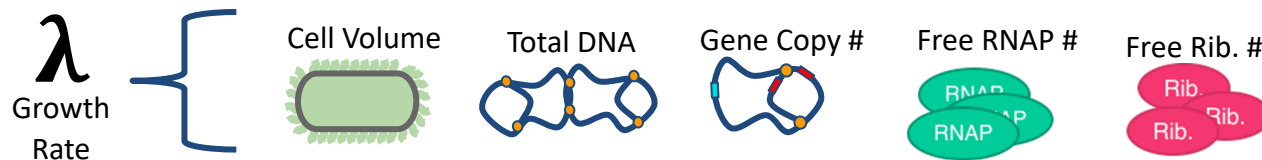
Dynamic stimulus methods

Optimal design: unravelling the effects of physiology on gene expression dynamics

Modelling goal: assess the effect of environmental factors on the dynamics of gene regulatory networks

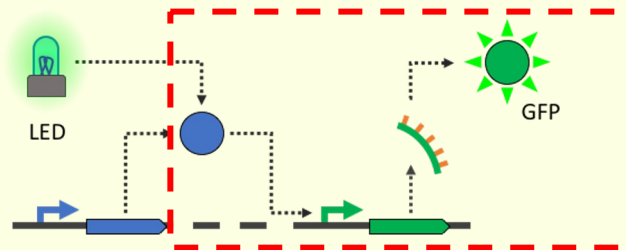


Klumpp and Hwa, (2014): Growth rate is a **sufficient statistic** for *E. coli* host physiology



Braniff, Scott, and Ingalls. Component characterization in a growth-dependent physiological context: optimal experimental design. *Processes* (2019)

Physiologically-aware gene expression model



Gene Expression Model

$$\frac{d X_{rna}}{dt} \frac{1}{V} = \alpha \frac{g}{V} \frac{\frac{P_a}{\eta G} K_r + \frac{P_a K_{rt}}{(\eta G)^2} u(t)}{1 + \frac{P_a}{\eta G} K_r + \left(\frac{K_i}{\eta G} + \frac{P_a K_{rt}}{(\eta G)^2} \right) u(t)} - \delta \frac{X_{rna}}{V}$$

$$\frac{d X_{prot}}{dt} \frac{1}{V} = \frac{\beta \frac{R_f}{V} X_{rna}}{K_M + \frac{R_f}{V}} - \lambda \frac{X_{prot}}{V}$$

Physiological Model Details

$$V = V_0 e^{(C+D)\lambda}$$

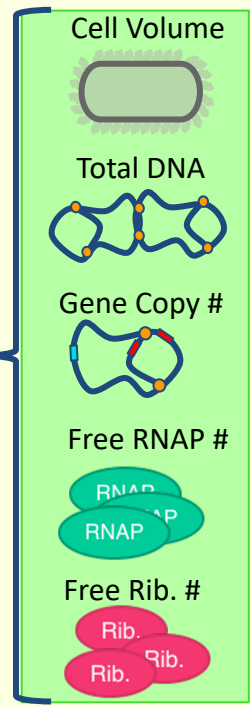
$$g = e^{((C+D) - l_{ori} C)\lambda}$$

$$G = \frac{1}{\lambda C} (e^{(C+D)\lambda} - e^{D\lambda})$$

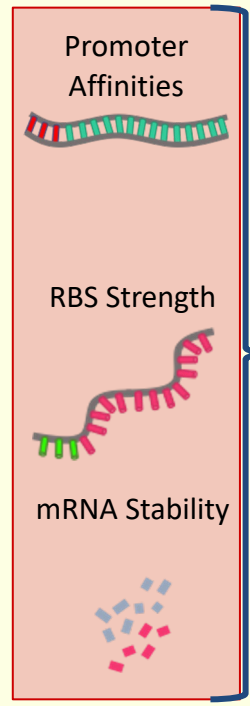
$$R_f = \frac{\rho V_0}{m_{rib}} \Phi_f (\kappa_r \lambda + \Phi_{r0}) (\kappa_{pr} \lambda + \Phi_{pr0}) e^{(C+D)\lambda}$$

$$P_a = \frac{\rho V_0}{m_{rna}} (\kappa_a \lambda + \Phi_{a0}) (\kappa_p \lambda + \Phi_{p0}) (\kappa_{pr} \lambda + \Phi_{pr0}) e^{(C+D)\lambda}$$

λ
Growth Rate



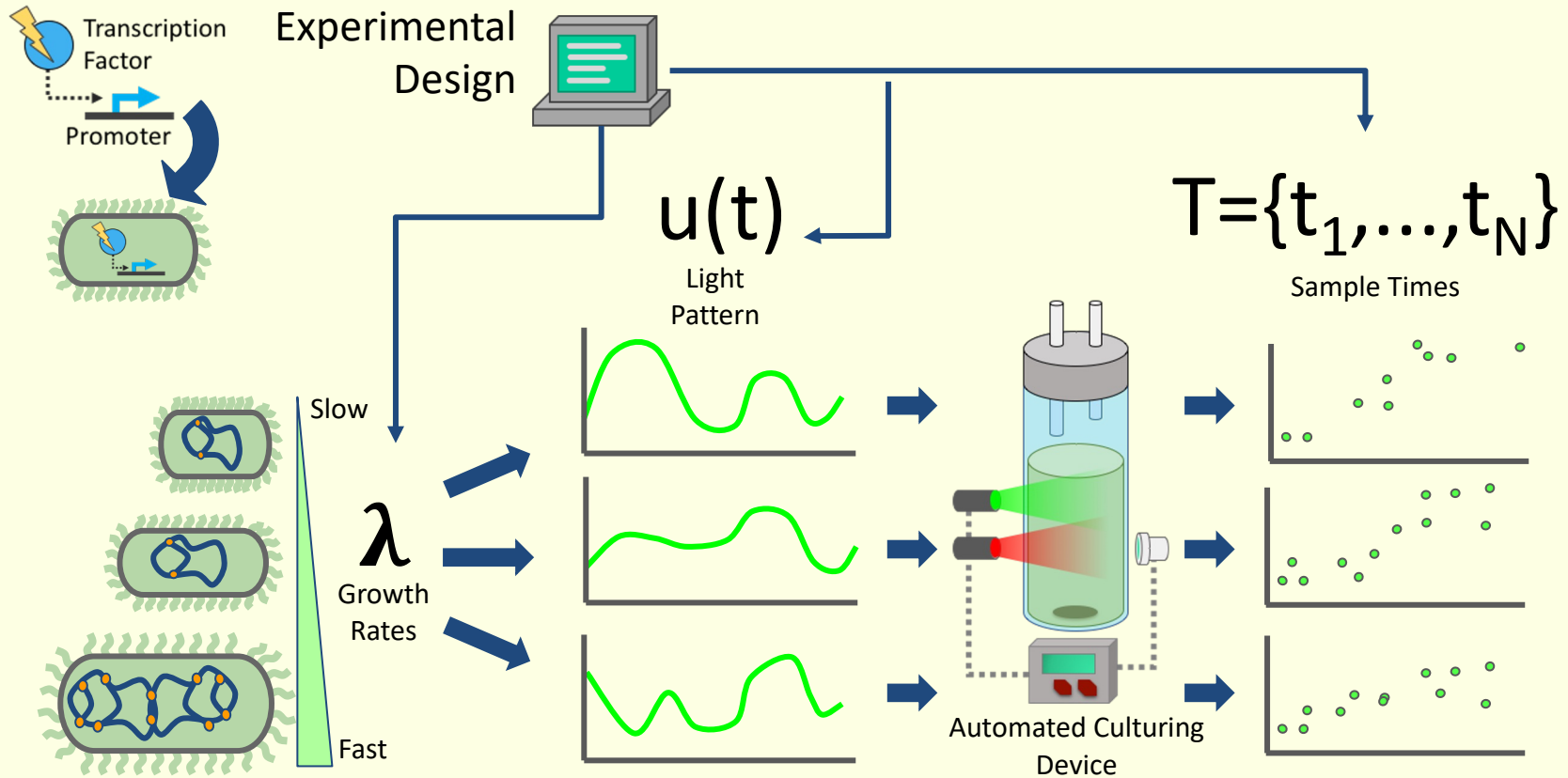
physiological parameters



gene-specific parameters

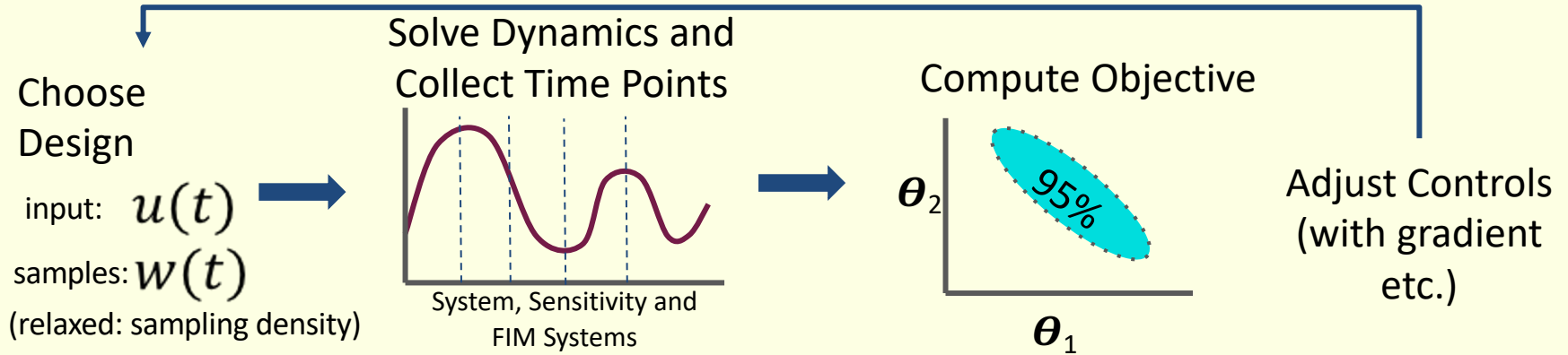
Intrinsic Parameters of the Sequence

Experimental Design

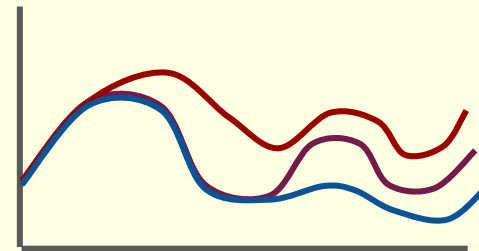
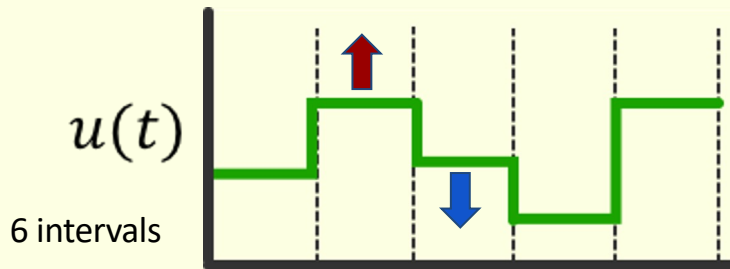


Optimization

Naive Approach



Computational challenges



Each control's effect is:

- 1) Distributed
- 2) Non-linear

Problem is poorly structured, gradient isn't useful

Solution: recast as optimal control problem

Experimental Design as Optimal Control

System Dynamics

$$\dot{\mathbf{y}} = \mathbf{F}(\mathbf{y}, \boldsymbol{\theta}, u(t), \lambda) \left\{ \begin{array}{l} \frac{d}{dt} \frac{X_{rna}}{V} = \alpha \frac{g}{V} \frac{\frac{P_a K_r}{\eta G} + \frac{P_a K_{rt} u}{(\eta G)^2}}{1 + \frac{P_a K_r}{\eta G} + \left(\frac{K_t}{\eta G} + \frac{P_a K_{rt}}{(\eta G)^2} \right) u} - \delta \frac{\xi}{V} \frac{X_{rna}}{V} \\ \frac{d}{dt} \frac{X_{prot}}{V} = \beta \frac{R_f}{V} \frac{X_{rna}}{V} - \lambda \frac{X_{prot}}{V} \end{array} \right.$$

$$\dot{\mathbf{S}} = \frac{\partial \mathbf{F}(\mathbf{y}, \boldsymbol{\theta}, u(t), \lambda)}{\partial \boldsymbol{\theta}} + \frac{\partial \mathbf{F}(\mathbf{y}, \boldsymbol{\theta}, u(t), \lambda)}{\partial \mathbf{y}} \mathbf{S} \quad \left(\mathbf{S} = \left[\frac{\partial \mathbf{y}}{\partial \theta_1}, \dots, \frac{\partial \mathbf{y}}{\partial \theta_N} \right] \right)$$

$$\dot{\mathbf{j}} = \mathbf{w}(t) \mathbf{S}^T \mathbf{S}$$

Telen et al., *Computers and Chemical Engineering*, 2014

Objective Function

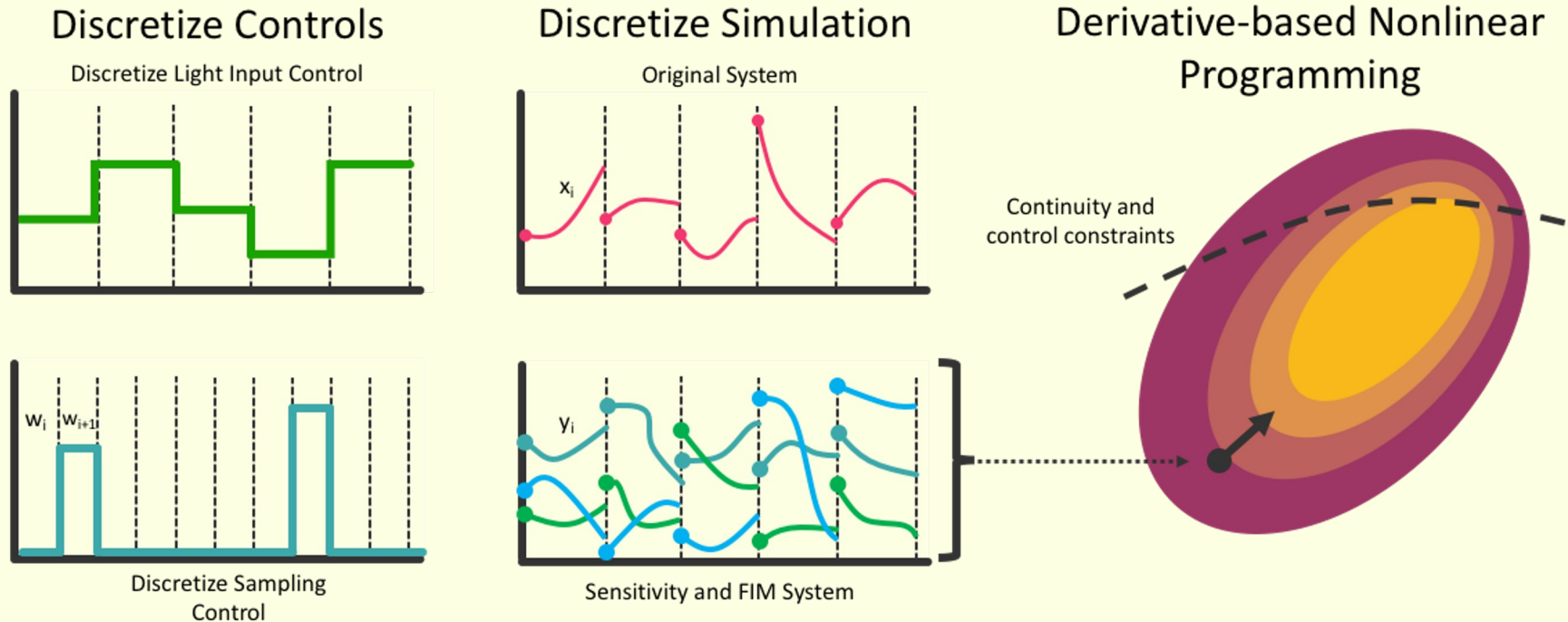
$$\det(\mathcal{J}(t_f))$$

Controls

$u(t)$ Induction Control λ Growth Rate

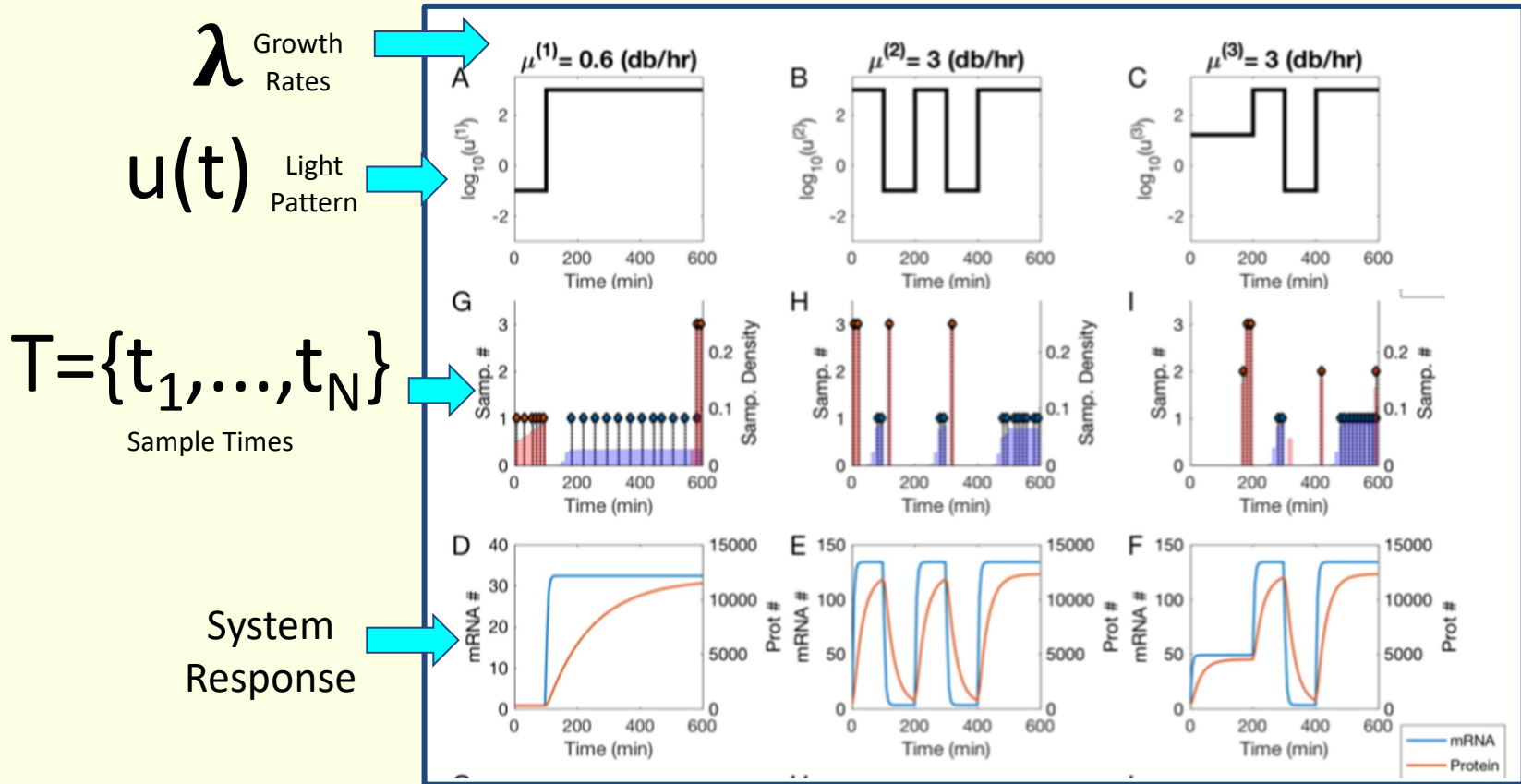
$w(t)$ Sampling (Continuous Relaxation)

Solution method: Multiple Shooting



CasADi symbolics for sensitivities of constraints (including initial conditions) and simulation steps (4th order Runge-Kutta)

Optimal Experiment

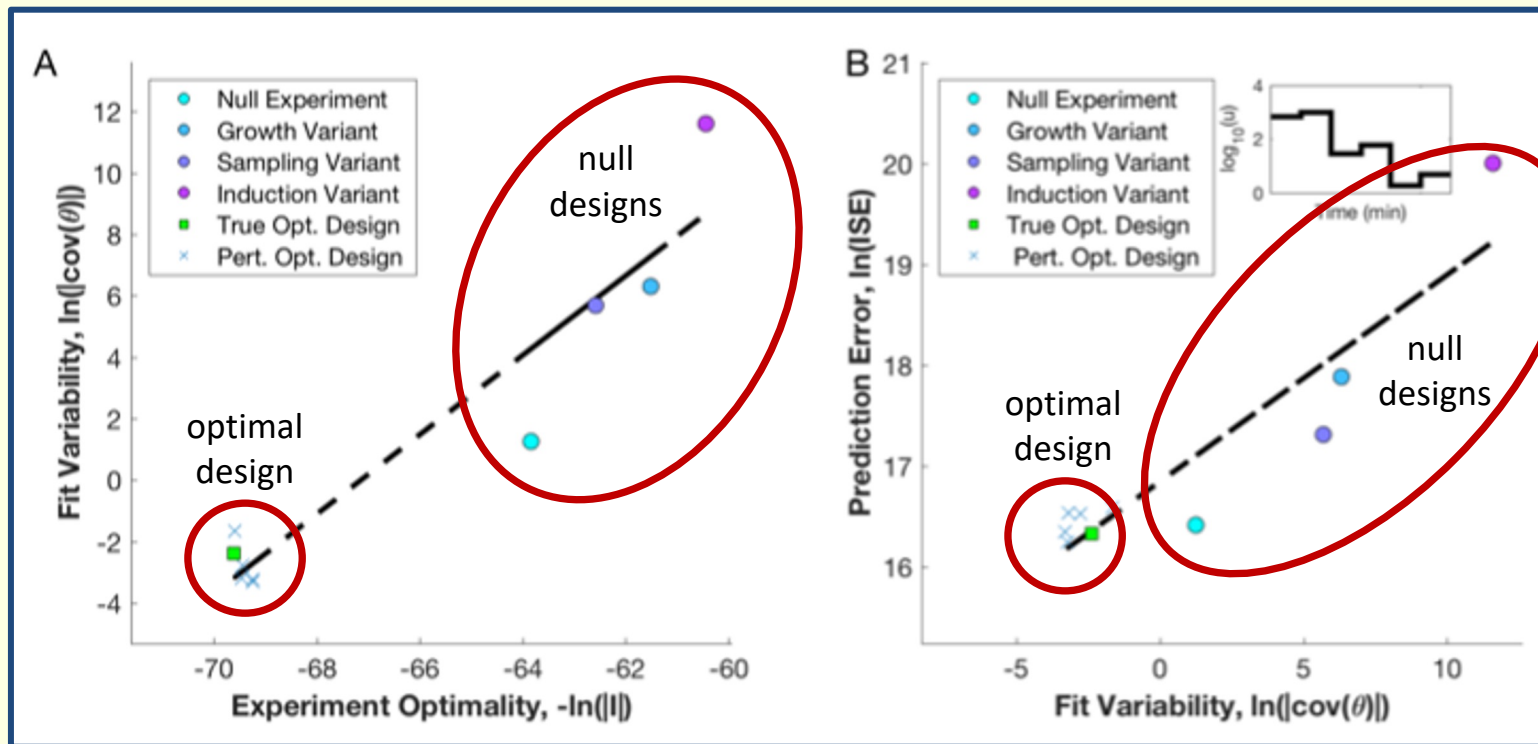


Validation

Parameter estimation

Prediction Accuracy
(out of sample experiment)

Parameter Variance



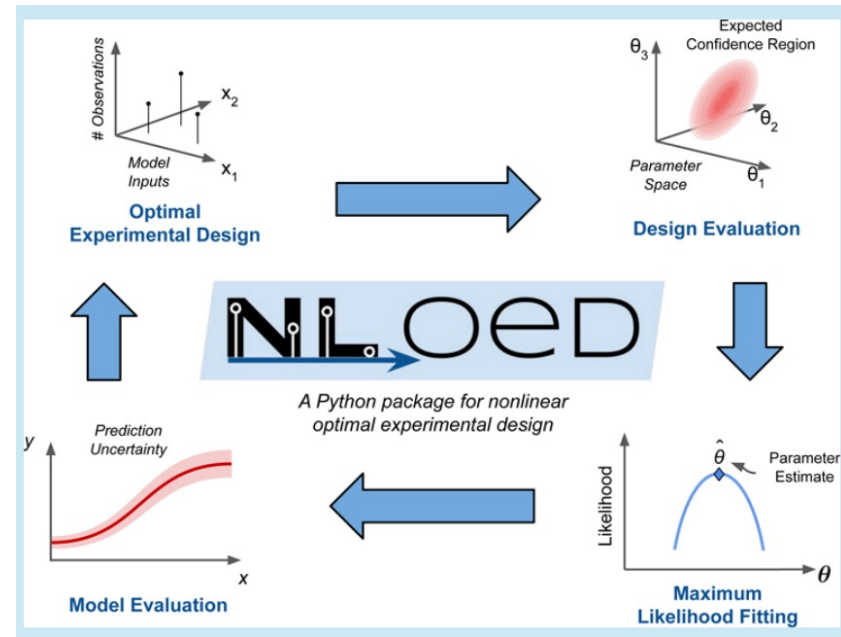
D-optimality score

Parameter Variance

Software package

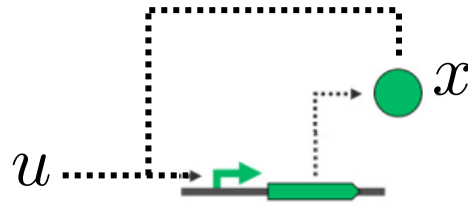
Python package: one-stop-shop for
FIM-based (local) Model-Based
Optimal Experimental Design

- Sequential design workflow
- Nonlinear models
- Non-Gaussian distributions, Poisson (e.g. plate counts), log-normal (e.g. gene expression), Bernoulli or binomial (e.g. viability assays)
- Symbolic model construction
- Sensitivities: automatic differentiation with CasADi
- Nonlinear programming: IPOPT
- D-optimal design over sampling and input profiles
- Integer sample counts relaxed to real-valued weights, then rounded
- Auxiliary methods: Maximum likelihood model fitting, sensitivity analysis, model simulation, data sampling, and design evaluation



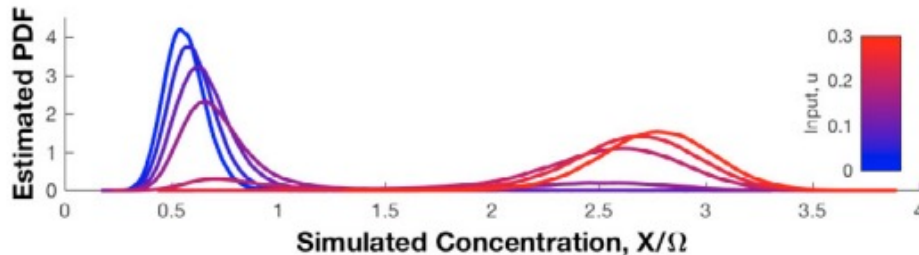
OED for multimodal gene expression system

Bistable autoactivating gene expression with activating external input

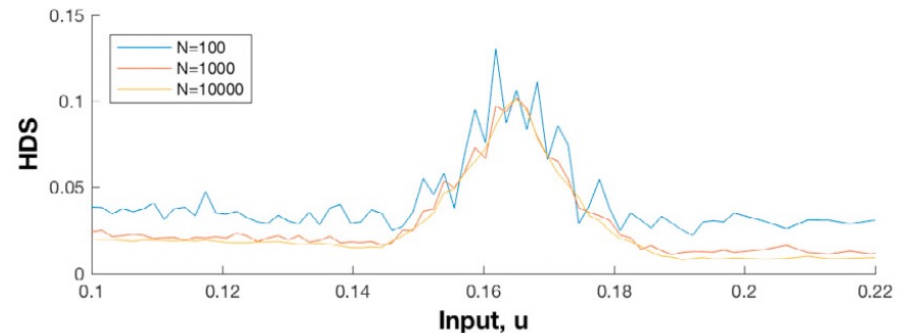


$$\frac{d}{dt}x(t) = \alpha_0 + \alpha \frac{(u + x(t))^n}{K^n + (u + x(t))^n} - x(t)$$

Observations bimodally distributed



Identified from steady state data via Hartigan dip statistic



Approximate log-likelihood: Gaussian mixture mode

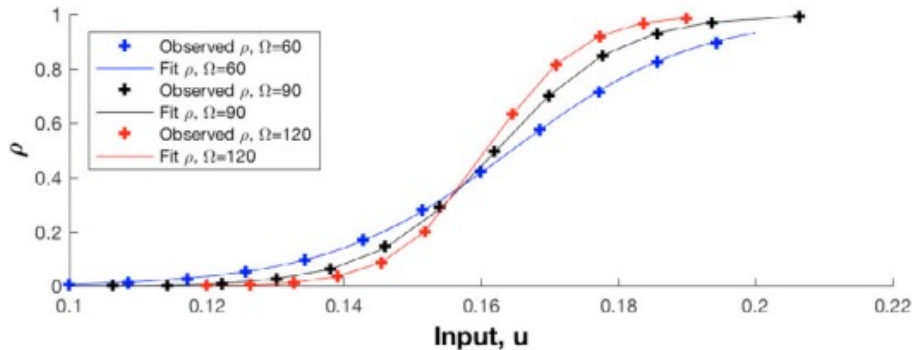
$$\ell(\theta|D, U) = \sum_i \log\{\rho(u_i) \cdot \varphi_T(y_i|u_i, \theta) + [1 - \rho(u_i)] \cdot \varphi_B(y_i|u_i, \theta)\}$$

OED for multimodal gene expression system

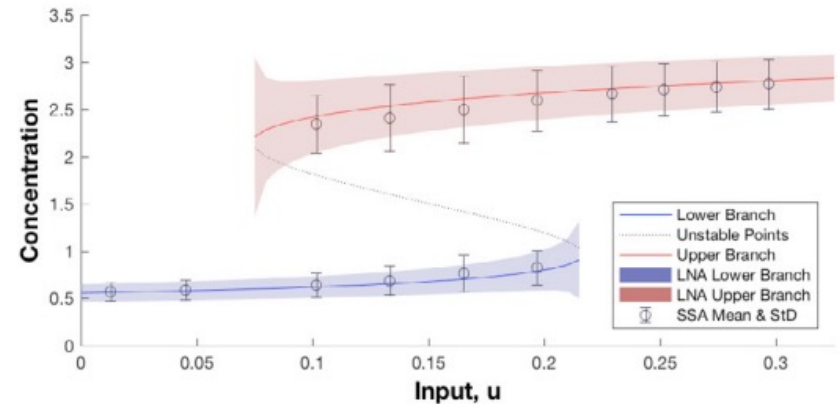
Approximate log-likelihood: Gaussian mixture mode

$$\ell(\theta|D, U) = \sum_i \log\{\rho(u_i) \cdot \varphi_T(y_i|u_i, \theta) + [1 - \rho(u_i)] \cdot \varphi_B(y_i|u_i, \theta)\}$$

Logistic approximation of probability of each mode (based on Kramers-Moyal approximation of escape times)



Linear Noise Approximation to estimate normal distribution around each mode:

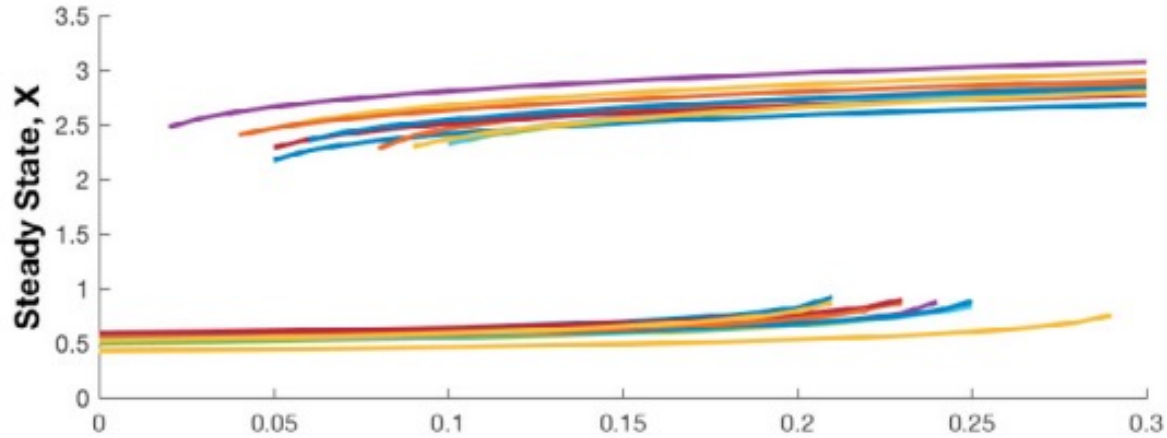


Validated by SSA

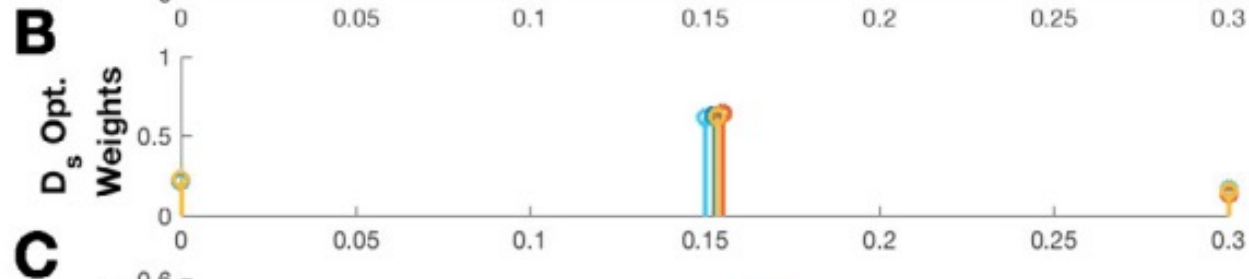
OED for multimodal gene expression system

Optimal designs

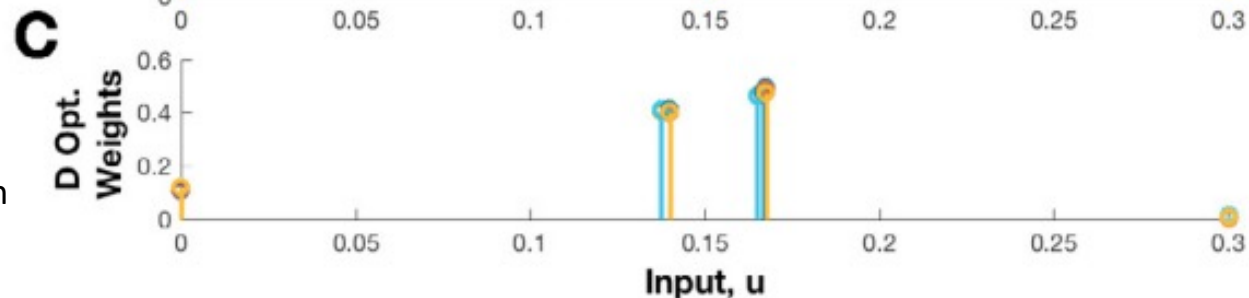
A



B



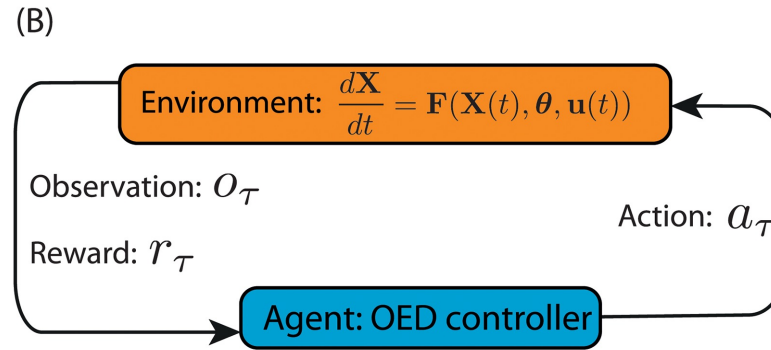
C



Deep reinforcement learning for OED

Reinforcement learning:

Agent receives **observation** and **reward**;
Implements **action** on **environment**



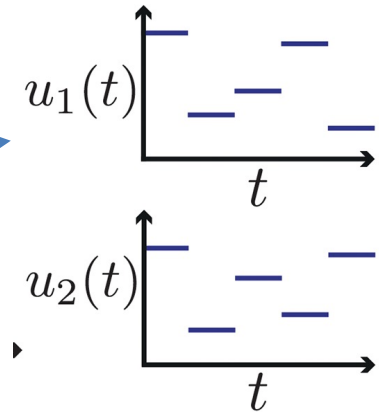
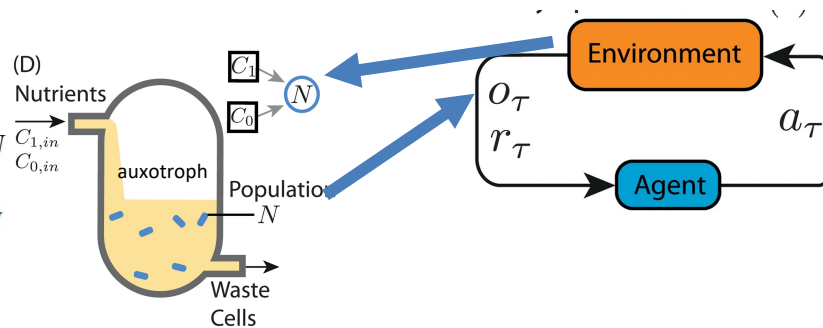
environment (system)
observation (time step, state measurement, estimate of FIM)
reward (increment in FIM)
action (experimental input)

$$\mu = \mu_{\max} \frac{C_1}{K_1 + C_1} \frac{C_0}{K_0 + C_0}$$

$$\frac{d}{dt} C_0 = q(C_{0,in} - C_0) - \frac{1}{\gamma_0} \mu N$$

$$\frac{d}{dt} C_1 = q(C_{1,in} - C_1) - \frac{1}{\gamma_1} \mu N$$

$$\frac{d}{dt} N = (\mu - q)N$$



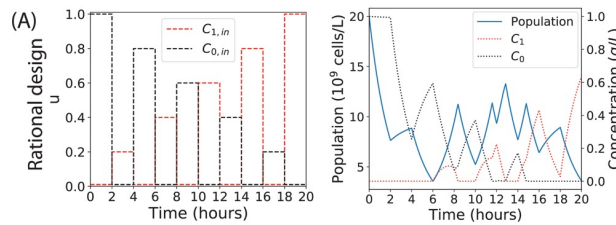
Advantage: learn model parameters while optimizing designs

Limitation: data hungry

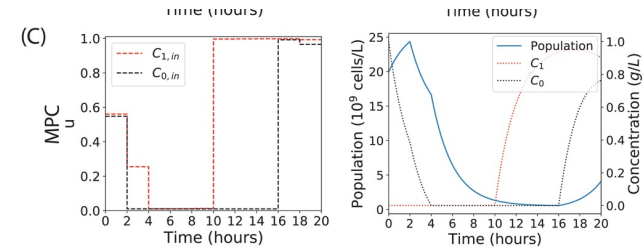
Fitted-Q learning over discrete action space: value function as deep neural network

Baseline performance: access to true parameter values

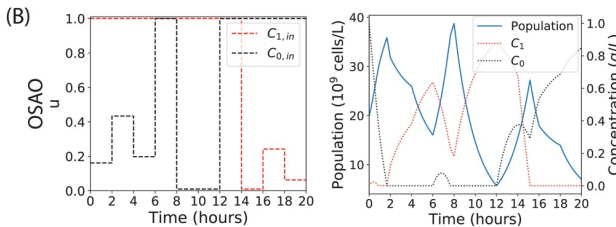
Rational (human) design



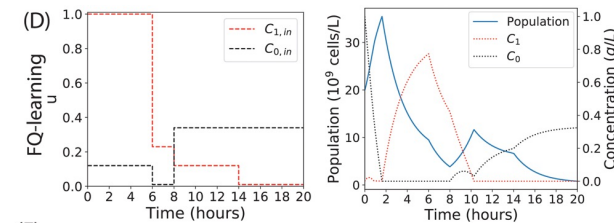
Model predictive control



Greedy optimization

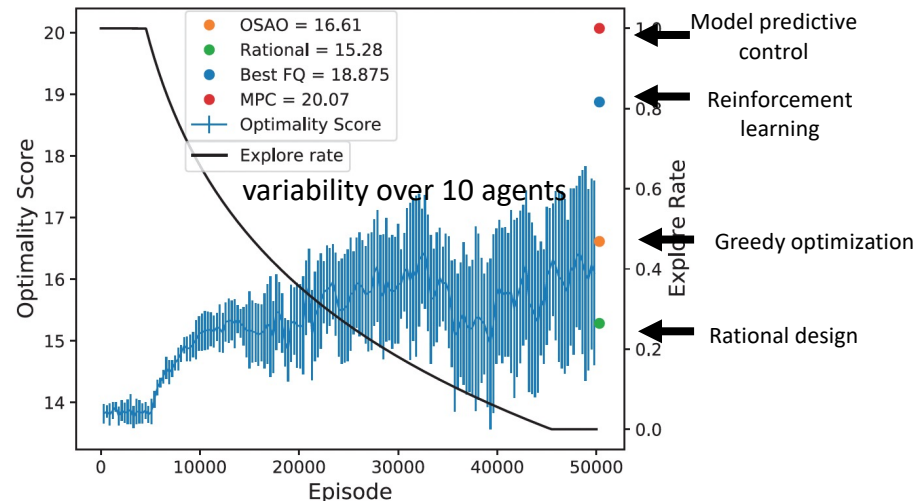


Reinforcement learning



Performance:

Training over 50000 simulated experiments

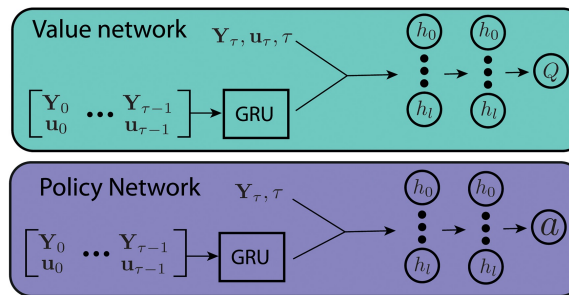
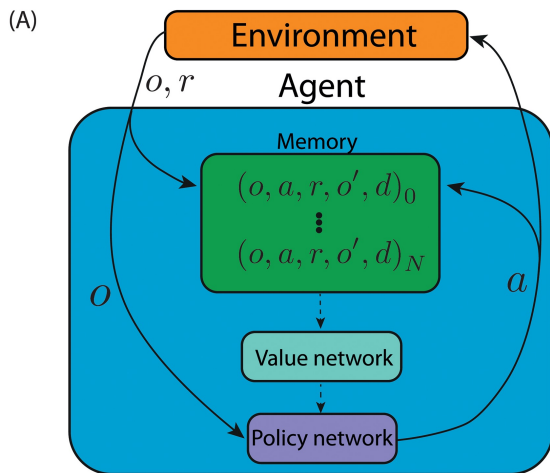


Similar results for parameter covariance and estimation accuracy

RL: promising results

Algorithm refinement: Recurrent Twin Delayed Deep Deterministic Policy Gradient (RT3D)

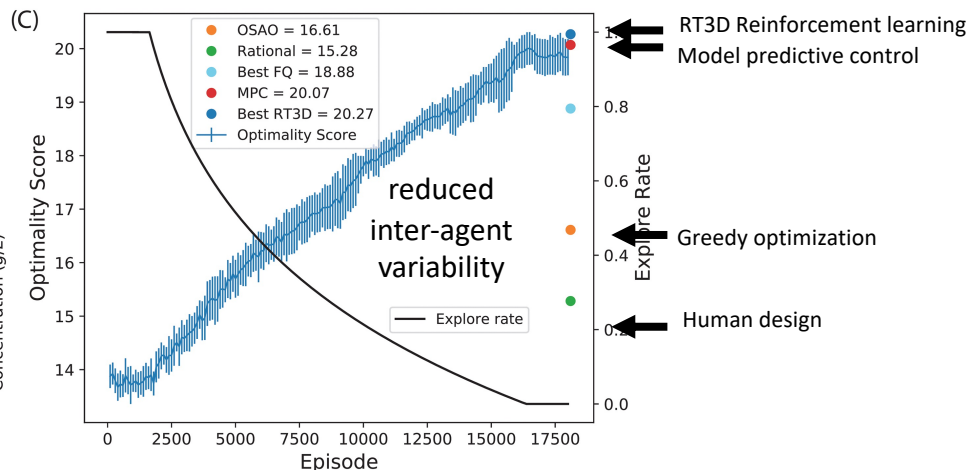
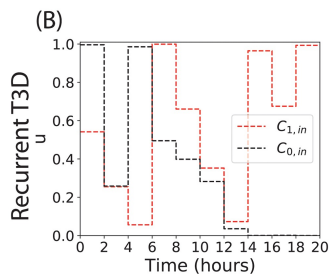
- Observation includes past history (allows learning of unknown parameter values)
- Continuous action space (requires additional recurrent neural network for feedback policy)



Baseline Performance

Training over 17500 simulated experiments

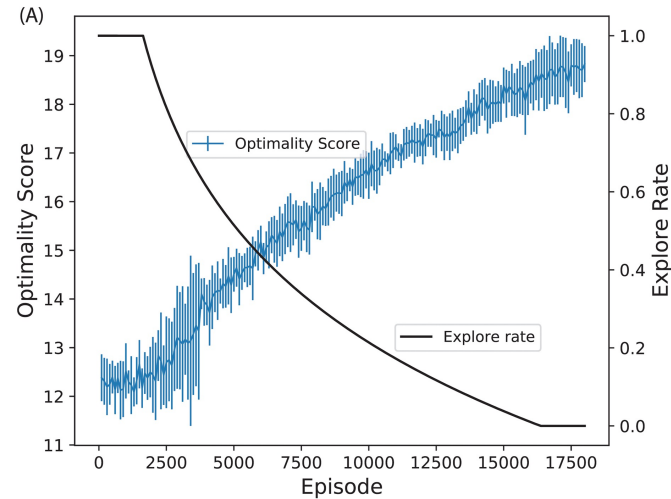
access to true parameter values



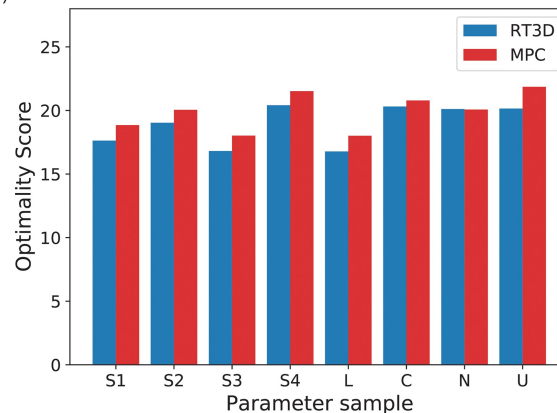
RL: equivalent to MPC

Agent performance over parameter distribution

10 agents. Each training simulation sampled from a uniform distribution



Performance comparison with **MPC acting with knowledge** of true sampled parameter values

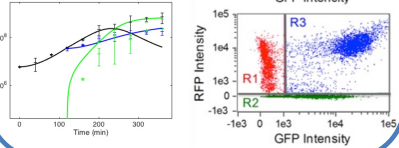


equivalent performance despite lack of a priori knowledge of parameter values

RL: improved robustness to parameter uncertainty in comparison to MPC

Acknowledgments

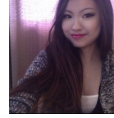
Nonspatial characterization



Akshay Malwade



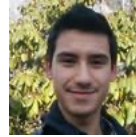
Angel Nguyen



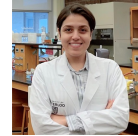
Peivand Sadat-Mousavi



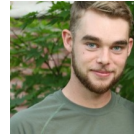
Aaron Yip



Atiyeh Ahmadi



Matt Courtney



Dhruva Rajwade



Nat Kendal-Freedman



Nathan Braniff



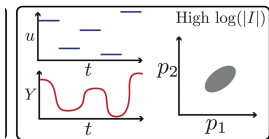
Addison Richards



Neythen Treloar



Experimental Design



Introductory modelling textbook:

PDF freely available at www.math.uwaterloo.ca/~bingalls/MMSB/

