# *Machine learning models applied to Synthetic Biology*

**Jean-Loup Faulon – Jean-Loup.Faulon@inrae.fr – https://jfaulon.com/course-materials**

**MICALIS Institute**

## Precise Prediction of Promoter Strength Based on a De Novo Synthetic Promoter Library Coupled with Machine Learning

Mei Zhao, Zhenqi Yuan, Longtao Wu, Shenghu Zhou*, and Yu Deng*

## Machine Learning Applied to Predicting Microorganism Growth Temperatures and Enzyme Catalytic Optima

Gang Li, Kersten S. Rabe, Jens Nielsen, and Martin K. M. Engqvist*

## SeqImprove: Machine-Learning-Assisted Curation of Genetic Circuit Sequence Information

Jeanet Mante, Zach Sents, Duncan Britt, William Mo, Chunxiao Liao, Ryan Greer, and Chris J. Myers*

## DNA Input Classification by a Riboregulator-Based Cell-Free Perceptron

Ardjan J. van der Linden, Pascal A. Pieters, Mart W. Bartelds, Bryan L. Nathalia, Peng Yin, Wilhelm T. S. Huck*, Jongmin Kim* and Tom F. A. de Greef*

## Design and Analysis of Compact DNA Strand Displacement Circuits for Analog Computation Using Autocatalytic Amplifiers

Tianqi Song, Sudhanshu Garg, Reem Mokhtar, Hieu Bui, and John Reif*

## Reservoir Computing Using DNA Oscillators

Xingyi Liu and Keshab K. Parhi*

## Reinforcement Learning for Bioretrosynthesis

Mathilde Koch, Thomas Duigou, and Jean-Loup Faulon*

## Generative Artificial Intelligence GPT-4 Accelerates Knowledge Mining and Machine Learning for Synthetic Biology

Zhengyang Xiao, Wenyu Li, Hannah Moon, Garrett W. Roell*, Yixin Chen*, and Yinjie J. Tang*

## DNA Memristors and Their Application to Reservoir Computing

Xingyi Liu and Keshab K. Parhi*

## Lessons from Two Design−Build−Test−Learn Cycles of Dodecanol Production in *Escherichia coli* Aided by Machine Learning

Paul Opgenorth, Zak Costello, Takuya Okada, Garima Goyal, Yan Chen, Jennifer Gin, Veronica Benites, Markus de Raad, Trent R. Northen, Kai Deng, Samuel Deutsch, Edward E. K. Baidoo, Christopher J. Petzold, Nathan J. Hillson, Hector Garcia Martin, and Harry R. Beller*

## Analog Computation by DNA Strand Displacement Circuits

Tianqi Song, Sudhanshu Garg, Reem Mokhtar, Hieu Bui, and John Reif*

## Machine Learning of Designed Translational Control Allows Predictive Pathway Optimization in *Escherichia coli*

Adrian J. Jervis, Pablo Carbonell, Maria Vinaixa, Mark S. Dunstan, Katherine A. Hollywood, Christopher J. Robinson, Nicholas J. W. Rattray, Cunyu Yan, Neil Swainston, Andrew Currin, Rehana Sung, Helen Toogood, Sandra Taylor, Jean-Loup Faulon, Rainer Breitling, Eriko Takano, and Nigel S. Scrutton*

## Tuning the Performance of Synthetic Riboswitches using Machine Learning

Ann-Christin Groher, Sven Jager, Christopher Schneider, Florian Groher, Kay Hamacher*, and Beatrix Suess*

## Semisupervised Gaussian Process for Automated Enzyme Search

Joseph Mellor, Ioana Grigoras, Pablo Carbonell, and Jean-Loup Faulon*

## Machine-Learning-Guided Mutagenesis for Directed Evolution of Fluorescent Proteins

Yutaka Saito, Misaki Oikawa, Hikaru Nakazawa, Teppei Niide, Tomoshi Kameda, Koji Tsuda*, and Mitsuo Umetsu*

## Supervised Learning in Adaptive DNA Strand Displacement Networks

Matthew R. Lakin* and Darko Stefanovic

# Supervised / Unsupervised Learning

**Precise Prediction of Promoter Strength Based on a De Novo Synthetic Promoter Library Coupled with Machine Learning**

Mei Zhao, Zhenqi Yuan, Longtao Wu, Shenghu Zhou*, and Yu Deng*

*ACS Synthetic Biology* 2022, 11, 1, 92-102 **(Research Article)**
Subscribed
**Publication Date (Web):** December 19, 2021
**DOI:** 10.1021/acssynbio.1c00117

**Semisupervised Gaussian Process for Automated Enzyme Search**

Joseph Mellor, Ioana Grigoras, Pablo Carbonell, and Jean-Loup Faulon*

*ACS Synthetic Biology* 2016, 5, 6, 518-528 **(Research Article)**
**Publication Date (Web):** March 23, 2016
**DOI:** 10.1021/acssynbio.5b00294

**Machine Learning Applied to Predicting Microorganism Growth Temperatures and Enzyme Catalytic Optima**

Gang Li, Kersten S. Rabe, Jens Nielsen, and Martin K. M. Engqvist*

*ACS Synthetic Biology* 2019, 8, 6, 1411-1420 **(Research Article)**

**Generative Artificial Intelligence GPT-4 Accelerates Knowledge Mining and Machine Learning for Synthetic Biology**

Zhengyang Xiao, Wenyu Li, Hannah Moon, Garrett W. Roell*, Yixin Chen*, and Yinjie J. Tang*

*ACS Synthetic Biology* 2023, 12, 10, 2973-2982 **(Research Article)**
Subscribed
**Publication Date (Web):** September 8, 2023
**DOI:** 10.1021/acssynbio.3c00310

**SeqImprove: Machine-Learning-Assisted Curation of Genetic Circuit Sequence Information**

Jeanet Mante, Zach Sents, Duncan Britt, William Mo, Chunxiao Liao, Ryan Greer, and Chris J. Myers*

*ACS Synthetic Biology*, **Articles ASAP (Technical Note)** Subscribed
**Publication Date (Web):** September 4, 2024
**DOI:** 10.1021/acssynbio.4c00392

**AlphaFold**

# Active / Reinforcement Learning

**Lessons from Two Design–Build–Test–Learn Cycles of Dodecanol Production in *Escherichia coli* Aided by Machine Learning**

Paul Opgenorth, Zak Costello, Takuya Okada, Garima Goyal, Yan Chen, Jennifer Gin, Veronica Benites, Markus de Raad, Trent R. Northen, Kai Deng, Samuel Deutsch, Edward E. K. Baidoo, Christopher J. Petzold, Nathan J. Hillson, Hector Garcia Martin, and Harry R. Beller*

*ACS Synthetic Biology* 2019, 8, 6, 1337-1351 **(Research Article)**

**Reinforcement Learning for Bioretrosynthesis**

Mathilde Koch, Thomas Duigou, and Jean-Loup Faulon*

*ACS Synthetic Biology* 2020, 9, 1, 157-168 **(Research Article)**
Subscribed
**Publication Date (Web):** December 16, 2019
**DOI:** 10.1021/acssynbio.9b00447

**Machine Learning of Designed Translational Control Allows Predictive Pathway Optimization in *Escherichia coli***

Adrian J. Jervis, Pablo Carbonell, Maria Vinaixa, Mark S. Dunstan, Katherine A. Hollywood, Christopher J. Robinson, Nicholas J. W. Rattray, Cunyu Yan, Neil Swainston, Andrew Currin, Rehana Sung, Helen Toogood, Sandra Taylor, Jean-Loup Faulon, Rainer Breitling, Eriko Takano, and Nigel S. Scrutton*

*ACS Synthetic Biology* 2019, 8, 1, 127-136 **(Research Article)**

**Tuning the Performance of Synthetic Riboswitches using Machine Learning**

Ann-Christin Groher, Sven Jager, Christopher Schneider, Florian Groher, Kay Hamacher*, and Beatrix Suess*

*ACS Synthetic Biology* 2019, 8, 1, 34-44 **(Research Article)**

**Machine-Learning-Guided Mutagenesis for Directed Evolution of Fluorescent Proteins**

Yutaka Saito, Misaki Oikawa, Hikaru Nakazawa, Teppei Niide, Tomoshi Kameda, Koji Tsuda*, and Mitsuo Umetsu*

*ACS Synthetic Biology* 2018, 7, 9, 2014-2022 **(Letter)**

**The Future of Go**

# in vitro / in vivo Learning

**Analog Computation by DNA Strand Displacement Circuits**

Tianqi Song, Sudhanshu Garg, Reem Mokhtar, Hieu Bui, and John Reif*

*ACS Synthetic Biology* 2016, 5, 8, 898-912 **(Research Article)**
**Publication Date (Web):** July 1, 2016
**DOI:** 10.1021/acssynbio.6b00144

**Design and Analysis of Compact DNA Strand Displacement Circuits for Analog Computation Using Autocatalytic Amplifiers**

Tianqi Song, Sudhanshu Garg, Reem Mokhtar, Hieu Bui, and John Reif*

*ACS Synthetic Biology* 2018, 7, 1, 46-53 **(Research Article)**
**Publication Date (Web):** December 4, 2017
**DOI:** 10.1021/acssynbio.6b00390

**DNA Input Classification by a Riboregulator-Based Cell-Free Perceptron**

Ardjan J. van der Linden, Pascal A. Pieters, Mart W. Bartelds, Bryan L. Nathalia, Peng Yin, Wilhelm T. S. Huck*, Jongmin Kim*, and Tom F. A. de Greef*

*ACS Synthetic Biology* 2022, 11, 4, 1510-1520 **(Research Article)** Open Access
**Publication Date (Web):** April 5, 2022
**DOI:** 10.1021/acssynbio.1c00596

**DNA Memristors and Their Application to Reservoir Computing**

Xingyi Liu and Keshab K. Parhi*

*ACS Synthetic Biology* 2022, 11, 6, 2202-2213 **(Research Article)** Subscribed
**Publication Date (Web):** May 13, 2022
**DOI:** 10.1021/acssynbio.2c00184

**Reservoir Computing Using DNA Oscillators**

Xingyi Liu and Keshab K. Parhi*

*ACS Synthetic Biology* 2022, 11, 2, 780-787 **(Research Article)** Subscribed
**Publication Date (Web):** January 26, 2022
**DOI:** 10.1021/acssynbio.1c00483

**Reference composition**
**Sun Z.Z. *et al. J. Vis. Exp.* 2013**

Lysate-based cell-free systems (TXTL)



Cell extract

- Grow and lyse cells
- Prepare crude extract
- Add substrates and salts
- Add template
- Incubate

| Component | Concentration | | | |
|---|---|---|---|---|
| Mg-glutamate (mM) | 0.4 | 1.2 | 2 | 4 |
| K-glutamate (mM) | 8 | 24 | 40 | 80 |
| Amino Acid (mM) | 0.15 | 0.45 | 0.75 | 1.5 |
| tRNA (mg/ml) | 0.02 | 0.06 | 0.1 | 0.2 |
| CoA (mM) | 0.026 | 0.078 | 0.13 | 0.26 |
| NAD (mM) | 0.033 | 0099 | 0.165 | 0.33 |
| cAMP (mM) | 0.075 | 0.225 | 0.375 | 0.75 |
| Folinic Acid (mM) | 0.0068 | 0.0204 | 0.034 | 0.068 |
| Spermidine (mM) | 0.1 | 0.3 | 0.5 | 1 |
| 3-PGA (mM) | 3 | 9 | 15 | 30 |
| NTP (mM) | 0.15 | 0.45 | 0.75 | 1.5 |

Combinatorial space = $4^{11}$
= **4 194 304** compositions

- Can we improve protein production without increasing the price of cell-free reaction?

- Can we provide efficient predictions of protein production *in vitro*?

- Can we highlight the critical parameters involve in protein production *in vitro*?

# *Active learning to optimize cell-free productivity*

- **Set up an initial batch sampling the space of possible compositions**
- **Measure yield level though fluorescence**
- **Develop a Neural Network models predicting yield from composition**
- **Use the models to predict the yield for each composition not yet tested**
- **Select next batch of compositions to be measured based on exploitation vs. exploration**
- **Repeat**





| | | |
|---|---|---|
| **Structural Materials Analysis** | Tools for High Energy X-ray Imaging, Diffraction, and Modeling of Microstructures | 6 |
| **Synthetic Biology** | Tools for synthetic biology | 64 |
| **Systems Biology** | Systems biology tools | 46 |

| | |
|---|---|
| icfree_instructor | iCFree instructor tool from the icfree suite |
| icfree_plates_designer | iCFree plates designer tool from the icfree suite |

- Herisson J, *et al. Nat Commun* **13**, 5082 (2022)

# Active learning to optimize cell-free productivity



• Borkowski O, *et al. Nat Commun* **11,** 1872 (2020)

6x more efficient than the best in vitro $CO_2$-fixing system described to date (CETCH 5.4 , Schwander *et al.* *Science* 2016)

- Pandi A., *et al. Nat Commun* **13,** 3876 (2022)

**TRAINING THE NETWORK**

Perceptron weights ($w_i$) are learned to increase classifier accuracy



*Input layer*

0.1  0.4  0.2 . .
.    .    .
.    .    .
0.3  0.0  0.1 . .
.    .    .
.    .    .
0.6  0.8  0.1 . .

$x_1$
$x_2$
$x_n$

*Weighted sum*
$w_1$
$w_2$
$w_n$

$b+\Sigma w_i x_i$

*Activation Function (sigmoid for classification)*
$f$
1
0

Prostate cancer

$$\text{PCa metabolic score} = b + \sum_{j=1}^{j} w_j x_{ij}$$



Ureido isobutyric acid
ROC Curve
Sensitivity = 70%
Specificity = 76.7%

Combined biomarker
ROC Curve
Sensitivity = 83.3%
Specificity = 83.3%

- Zang*, et al. PLoS One* 2013 and *J Proteome Res.* 2014
- Shen B, et al. *Cell*. 2020
- …..

**USING THE TRAINED NETWORK**

To perform a diagnostic:

- Quantify a panel of biomarkers (metabolites) on clinical samples (using metabolomics)
- Feed measured biomarkers concentrations ($x_i$) to

$f\ (b+\Sigma w_i x_i\ )$

- **Is it possible to avoid biomarker concentration measurements?**

  ➢ **Engineer the trained network *in vitro* or *in vivo* and directly use it on clinical samples**

# Engineering a neural metabolic network: the concept



## TRAINING THE NETWORK

Perceptron weights ($w_i$) are learned to increase classifier accuracy

Input layer

Weighted sum

Activation Function (sigmoid for classification)

$x_1$  $w_1$  $x_2$  $w_2$  $b+\Sigma w_i x_i$  $f$  $x_n$  $w_n$

0.1  0.4  0.2 ..
0.3  0.0  0.1 ..
0.6  0.8  0.1 ..

1
0

Prostate cancer

$$\text{PCa metabolic score} = b + \sum_{j=1}^{j} w_j x_{ij}$$

Ureido isobutyric acid
ROC Curve

Sensitivity = 70%
Specificity = 76.7%

Combined biomarker
ROC Curve

Sensitivity = 83.3%
Specificity = 83.3%

## ENGINEERING THE TRAINED NETWORK

Need to actuate weighted sum and activation function

$x_1$  $E_1$  Activated-TF  GFP  $x_2$  $E_2$  Activator  $b+\Sigma w_i x_i$  TF  $E_n$  $f(b+\Sigma w_i x_i)$  $x_n$

Enzymatic transformation     Reporter gene

In theory (Michaelis-Menten) when $x_i << [E_i]$ :

$$d[p] = \Sigma k_i [E_i] x_i \, dt$$

$w_i = k_i [E_i]$
where $k_i$ is a kinetics constant

Sigmoid behavior

Fold Change
[Benzoic Acid] (µM)

- Zang, *et al.* **PLoS One** 2013 and ***J Proteome Res.*** 2014
- Shen B, et al. **Cell**. 2020
- .....

**ENGINEERING THE TRAINED NETWORK**

Need to actuate weighted sum and activation function

$x_1$

$E_1$

$E_2$

Activated-TF

GFP

$x_2$

Activator

$b + \Sigma w_i x_i$

TF

$E_n$

$f(b + \Sigma w_i x_i)$

$x_n$

*Enzymatic transformation*     *Reporter gene*

Undetectable Inducer

Benzoate     Benzoate     sfGFP

OR2-OR1-Pr Enz     OR2-OR1-Pr BenR     BenR     sfGFP     BenR sfGFP
TF plasmid     P_Ben Reporter plasmid     P_Ben

$$\frac{dbenzoate}{dt} = enz * \frac{k_{cat} * inducer}{inducer + K_M}$$

$$\frac{dinducer}{dt} = -enz * \frac{k_{cat} * inducer}{inducer + K_M}$$

$$TF_{activated} = TF * \frac{benzoate}{benzoate + K_d^{inducer}} + 0.0005$$

$$\epsilon = \frac{TF_{activated}}{TF_{activated} + K_d^{activated}} \text{ for BenR}$$

$$\epsilon = 1 \text{ for constitutive expression}$$

$$\frac{dmRNA}{dt} = \gamma * n * \epsilon \frac{x}{x + \chi} * \frac{K_{tox}}{K_{tox} + tox} * \frac{R_{mRNA}}{R_{mRNA} + K_{mRNA}} - \delta * mRNA$$

$$\frac{dprot}{dt} = \pi * mRNA * \frac{y}{y + k} * \frac{K_{tox}}{K_{tox} + tox} - \lambda * prot$$

measurement     model     $x_i$

$w_i$

[Hippuric acid] (µM)
1000 500 200 100 50 20 10 5 2 1 0.5 0

REU
40 35 30 25 20 15 10 5

0 0.1 0.3 1 3 10 30 100
[HipO DNA] (nM)

0 0.1 0.3 1 3 10 30 100
[HipO DNA] (nM)

In theory (Michaelis-Menten)
when $x_i << [E_i]$ :

$$d[p] = \Sigma k_i [E_i] x_i \, dt$$

$w_i = k_i [E_i]$
where $k_i$ is a kinetics constant

*Sigmoid behavior*

Fold Change
300 100 30 10 3 1

1 10 100 1000
[Benzoic Acid] (µM)

• Voyvodic, P.L., Pandi, A., Koch, M. *et al. Nat Commun* **10**, 1697 (2019).

# Engineering a neural metabolic network in vitro

TARGETED BEHAVIOUR

Hippurate
Cocaine
Benzamide
Biphenyl-2,3-diol

$f(b + \sum w_i x_i) < 0.5$    $f(b + \sum w_i x_i) \geq 0.5$

- Kinetics model is used to compute the enzyme concentration for each weight

$f(b + \sum w_i x_i)$
$b = -0.50$
$w_1 = 0.25$
$w_2 = 0.25$
$w_3 = 0.50$
$w_4 = 0.50$

Logistic regression

Hippurate
Cocaine
Benzamide
Biphenyl-2,3-diol

$E_1$  $E_2$  $E_3$  $E_{42}$  $E_{41}$  Benzoate  BenR  Activated-TF  GFP

Retrosynthesis Workflow »    Galaxy

[Enzyme DNA] nM
[Hippurate] µM
[Cocaine] µM
[Benzamide] µM
[Biphenyl-2,3-diol] µM
REU

TARGETED BEHAVIOUR

(−)    (+)

RFU

(+)
(−)

OBSERVED BEHAVIOUR

Hippurate
Cocaine
Benzamide
Biphenyl-2,3-diol

- Pandi A., Koch M. *et al*. *Nat Commun* **10**, 3880 (2019)

# Engineering a neural metabolic network in vitro



TARGETED BEHAVIOUR

Hippurate..............
Cocaine...................
Benzamide..............
Biphenyl-2,3-diol...

$f(b + \sum_i w_i x_i) < 0.5$     $f(b + \sum_i w_i x_i) \geq 0.5$

- Kinetics model is used to compute the enzyme concentration for each weight

$f(b + \sum_i w_i x_i)$
$b = -0.50$

Logistic regression

$w_1 = \cancel{0.25}\, 0.50$
$w_2 = \cancel{0.25}\, 0.50$
$w_3 = \cancel{0.50}\, 0.50$
$w_4 = \cancel{0.50}\, 0.50$

Hippurate
Cocaine
Benzamide
Biphenyl-2,3-diol

Activated-TF   GFP

$E_1$
$E_2$
$E_3$
$E_{42}$
$E_{41}$
Benzoate   BenR

Retrosynthesis Workflow »   Galaxy

TARGETED BEHAVIOUR

OBSERVED BEHAVIOUR

Hippurate..............
Cocaine...................
Benzamide..............
Biphenyl-2,3-diol...

- Pandi A., Koch M. *et al.* *Nat Commun* **10**, 3880 (2019)

# Engineering a neural metabolic network in vivo?



- Can we divert native metabolism to handle problems that are usually solved *in silico*?

- Ability of physical, chemical or biological devices to solve problems is studied in AI with Reservoir Computing (RC)



(a) Conventional RC

(b) Physical RC

Tanaka G. et al. *Neural Networks* **115**, 100 (2019)

# *E. coli Reservoir Computer (E. coli RC)*

Can we exploit *E. coli* native metabolism to build an *E. coli* RC to solve computational problems?

How complex a problem can *E. coli* RC solve?

Can we find practical uses of *E. coli* RC?

# E. coli Reservoir Computer (E. coli RC)

Can we exploit *E. coli* native metabolism to build an *E. coli* RC to solve computational problems?

How complex a problem can *E. coli* RC solve?

Can we find practical uses of *E. coli* RC?



gradient backpropagation

Conventional Reservoir should:
- accurately reproduce phenotype for different media composition
- enable gradient backpropagation

# The reservoir

GEnome-scale Metabolic Model (GEM/FBA)

Max ($v_{biomass}$)

Subjected to contraints:
    S $V = 0$
    $0 \leq V \leq V_{in}$

where
$- V$ = set of all reaction fluxes
$- S$ = stochiometric matrix
$- V_{in}$ = uptake medium fluxes upper bounds



Conventional Reservoir

# *The reservoir*

GEnome-scale Metabolic Model (GEM/FBA)

Max ($v_{biomass}$)

Subjected to contraints:
   $S\,V = 0$
   $0 \leq V \leq V_{in}$

where
$- V$ = set of all reaction fluxes
$- S$ = stochiometric matrix
$-V_{in}$ = uptake medium fluxes upper bounds

*GEM/FBA growth rates vs. measured growth rate in E. coli MG1655 for 1 to 4 nutrients added to M9*



Conventional Reservoir



Medium concentrations

?

Medium fluxes

Reaction fluxes

Growth rate

*Conventional Reservoir should:*
- *accurately reproduce phenotype for different media composition*

# Building an E. coli RC to increase mechanistic model predictability

GEnome-scale Metabolic Model (GEM/FBA)

$\text{Max } (v_{biomass})$

Subjected to contraints:
$$S\,V = 0$$
$$0 \le V \le V_{in}$$

where
- $V$ = set of all reaction fluxes
- $S$ = stochiometric matrix
- $V_{in}$ = uptake medium fluxes upper bounds

*GEM/FBA is a Linear Program solved using Simplex algorithm not compatible with gradient propagation*

*GEM/FBA growth rates vs. measured growth rate in E. coli MG1655 for 1 to 4 nutrients added to M9*



**Prior-ANN   Physical  Reservoir   Post-ANN**

Conventional Reservoir



*gradient backpropagation to find mapping between medium concentrations and uptake fluxes*

*Conventional Reservoir should:*
- *accurately reproduce phenotype for different media composition*
- *enable gradient backpropagation*

## Physics informed neural network (PINN)



$V_{in}$

$V^0 \leftarrow W_i V_{in}$ — **Neural layer**

$V \leftarrow V + \nabla V$ — **Mechanistic layer**

$\nabla V$ is computed from $S$

$V_{out}$

## Hopfield's network



$V_{in}$

$V^0 \leftarrow W_i V_{in}$ — **Neural layer**

$U, V$

$U \leftarrow U + \nabla U$
$V \leftarrow V + \nabla V$ — **Mechanistic layer**

$\nabla V$ and $\nabla U$ are computed from $S$

$U^0 = 0$
$U$ are variables of the dual problem (shadow metabolites)

$V_{out}$

## Graph neural network (GNN)



$V_{in}$

$V^0 \leftarrow W_i V_{in}$ — **Neural layer**

$V$

$M \leftarrow P_{v \to m} V$
$V \leftarrow (P_{m \to v} \odot W_r) M + V^0$ — **Mechanistic layer**

$W_r$ are flux split ratios at branch metabolites

$P_{v \to m}$ and $P_{m \to v}$ are adjacency matrices computed from $S$

$V_{out}$

Trained on GEM/FBA ( cobrapy ) calculated growth rates with *E. coli* iML1515 model for 1000 different media (M9 + random combinations of nutrients among sugars, nucleotides, amino acids)







- Faure L. *et al. Nat Commun* **14**, 4669 (2023)

Medium = M9 + 280 combinations of nutrients with fixed concentration (among 28 sugars, nucleotides, amino acids)

medium composition

Ala    Gua    Glc

Scaler to medium fluxes

cobrapy
*E. coli* GEM/FBA

Calculated growth rate

$R^2 = 0.12$

Calculated growth rate (h$^{-1}$)

Measured growth rate (h$^{-1}$)

**GEM/FBA results with best scaled input**

- Faure L. *et al*. *Nat Commun* **14**, 4669 (2023) & Faulon et al. *bioRxiv DOI: 10.1101/2024.09.12.612674* (2024)

GEM/FBA results with best scaled input

$R^2 = 0.12$

GEM/FBA results with reservoir inputs

$R^2 = 0.98$

Medium = M9 + 280 combinations of nutrients with fixed concentration (among 28 sugars, nucleotides, amino acids)

Training on 280 medium compositions

Reservoir
Trained on GEM/FBA calculated growth rates

$V_{in}$

$V_{out}$

Predicted growth rate

Calculated growth rate

- Faure L. *et al*. *Nat Commun* **14**, 4669 (2023) & Faulon et al. *bioRxiv DOI: 10.1101/2024.09.12.612674* (2024)

Medium = M9 + 280 combinations of nutrients with fixed concentration (among 28 sugars, nucleotides, amino acids)

**Training-set 80% - Test-set 20% (5-fold CV)**

medium composition

**Conventional RC**

**Neural layer**

$V_{in}$

**Test-set**

**Reservoir**
*Trained on GEM/FBA calculated growth rates*

**Neural layer**

**AMN**

**Mechanistic layer**

$V_{out}$

Predicted growth rate

~

**cobrapy**
*E. coli* **GEM/FBA**

Calculated growth rate

**Test-set**

**R² = 0.12**

GEM/FBA results with best scaled input

**R² = 0.78**

GEM/FBA results with reservoir inputs

• Faure L. *et al*. *Nat Commun* **14**, 4669 (2023) & Faulon et al. *bioRxiv DOI: 10.1101/2024.09.12.612674* (2024)

# Can E. coli RC be used to solve a classical machine learning problem?

# Using E. coli RC to solve a regression problem

**Example of regression problem : OpenML 'Energy Efficiency' dataset (768 instances, X = 8 features, y = % efficiency)**



| Compactness | Surface Area | Orientation |
|---|---|---|
| 0.98 | 514.5 | 3 |

| Glc | Xyl | Succ | Trp |
|---|---|---|---|
| 1 | 0 | 1 | 1 |

| μ |
|---|
| 0.71 |

| % efficiency |
|---|
| 15.55 |

• Faulon et al. *bioRxiv DOI: 10.1101/2024.09.12.612674* (2024)

# Using E. coli RC to solve a regression problem

**Example of regression problem : OpenML 'Energy Efficiency' dataset (768 instances, X = 8 features, y = % efficiency)**



- Faulon et al. *bioRxiv DOI: 10.1101/2024.09.12.612674* (2024)

OpenML
A worldwide machine learning lab

**10 OpenML regression problems of increasing difficulty**



- Faulon et al. *bioRxiv DOI: 10.1101/2024.09.12.612674* (2024)

**10 OpenML regression problems of increasing difficulty**

Legend:
- MLR
- MLP
- XGB
- Conventional RC (generative prior)
- Physical RC (selective prior)

• Faulon et al. *bioRxiv DOI: 10.1101/2024.09.12.612674* (2024)

• Faulon et al. *bioRxiv DOI: 10.1101/2024.09.12.612674* (2024)

• All datasets from Baltussen et al. Nature 2024

The problem:

- Blood sample are collected for Covid-19 patients once they enter the hospital

- Metabolomics analyses are carried out on the samples

- Can we predict from the analyses if the disease outcome will be severe or moderate?

CHU Grenoble-Alpes cohort (training set):

- 81 patients
- 624 molecules detected (56 *E. coli* medium molecules)
- severe (31) – moderate (50)

### Classifier performances (20-fold CV results)



Accuracy = 0.84 in Shen et al. *Cell* 2020; 182(1): 59–72

# Using E. coli RC for classification

The problem:

- Blood sample are collected for Covid-19 patients once they enter the hospital

- ~~Metabolomics analyses are carried out on the samples~~

- ~~Can we predict from the analyses if the disease outcome will be severe or moderate?~~

- Can we use an *E. coli* RC grown on the patient's sample to predict if the disease outcome will be severe or moderate ?

Conventional RC to predict disease outcome from phenotype

Prior input     Reservoir     Post readout

*Metdium metabolite MS signals*

*E. coli* MG1655 or *E. coli* gene-KO AMN model (GEM iML1515)

*Reaction fluxes*

moderate vs. severe

*Medium metabolites uptake rates*

### MG1655 strain

True positive rate vs. False positive rate

AUC = 0.66

### Gene KO strains

Accuracy: MG1655, Δgene1, Δgene2, Δgene3, Δgene4, Δgene5, Δgene6

# Building an E. coli physical RC for classification



- ➢ Gene-KO *E. coli* Physical RC to predict disease outcome from growth rate and $OD_{MAX}$

- ➢ CRISPR-Cas9/Lambda red system

  - o Jiang et al. , *Appl Environ Microbiol*, 2015
  - o Scarless, Efficient, Multiplexable
  - o ~80 KOs built

- ➢ Gene deletions force *E. coli* to collect specific nutrients from the plasma in order to grow

- ➢ According to conventional RC, differences of nutrients concentration in the plasma should result in different growth curve

# *Perspectives*

- Supervised Learning & Active learning

  - o New generative AI models (transformer) for retro-(bio)synthesis and to generate sequences
  - o LLMs (like GPT4) to drive biofoundries
  - o Active Learning / Transfer Learning / Hybrid Learning to cope with small training set sizes

- *in vitro/in vivo* learning

  - o Decades of research and development in Synthetic Biology to build bottom-up computing devices (digital, analog, neural,…)… but many difficulties
  - o Most devices were inspired from natural biological networks: perhaps one should to consider building devices top-down, *i.e.* exploiting/modifying hosts rather than plugging orthogonal devices.

# *Acknowledgments*

ANR

AMN
SynBioDiag

B-BEST

FRANCE 2030

HORIZON    BIOS
SUSTAINABLE BIO-MANUFACTURING

université PARIS-SACLAY

INRAE

DGA    X
ÉCOLE POLYTECHNIQUE

Paul Ahvi

An Hoang

Bastien Mollet

Léon Faure

Amir Pandi

Mathilde Koch

Olivier Borkowski

Anne Giralt

Joan Hérisson

Jérôme Bonnet's group

CBS    Inserm

Tobby Erb's group

MAX PLANCK INSTITUTE
FOR TERRESTRIAL MICROBIOLOGY

⭐ Molecular Biology
⭐ Computational Biology