



13th International Workshop on Bio-Design Automation
Online
September 20th-24th 2021

Foreword

Welcome to IWBD A 2021!

The IWBD A 2021 Organizing Committee welcomes you to the Thirteenth International Workshop on Bio-Design Automation (IWBD A). IWBD A brings together researchers from an array of life science and computational disciplines including synthetic biology, systems biology, and design automation. The focus is on concepts, methodologies, and tools for the synthesis of biological systems.

While the field of synthetic biology is still nascent, we can already appreciate that many success stories can be attributed to the joint efforts of researchers with experimental expertise and researchers making computational contributions. As we work to convert the design of synthetic biological systems from an ad hoc process to one driven by standardized components and principles, there is a tremendous opportunity to engage with new ideas and new communities. IWBD A offers a forum for cross-disciplinary discussion, with the aim of seeding and fostering collaboration between biological and computational research communities. We hope that over the next few days, you will all encounter new ideas you can integrate into your current efforts.

This year, the program consists of 8 workshops, 13 contributed talks, and 7 short presentations (in lieu of posters): The talks are organized into 6 sessions:

- Standardization of Biological Components
- Screening Methods
- Metabolic and Knowledge Engineering
- Computer-aided design, modelling, and simulation
- Machine Learning
- Data Repositories

In addition, we are very pleased to have two distinguished invited speakers: Dr. Tijana Radivojevic from the Lawrence Berkeley National Laboratory and Dr. Thomas E. Goro-chowski from the University of Bristol

IWBD A is proudly organized by the non-profit Bio-Design Automation Consortium (BDAC). BDAC is an officially recognized 501(c)(3) tax-exempt organization.

We would like to thank all the participants for their contributions to IWBD A. We would also like to highlight the efforts of the Program Committee and the Virtual Chairs. Our combined efforts enable us to adapt to a more accessible online format for this year.

Organizing Committee

Organizing Committee

General Chair - Prashant Vaidyanathan, Microsoft Research

Program Committee Chair - Marilene Pavan, LanzaTech

Publication Chair - Jenhan Tao, Generate Biomedicines

Logistics Chair - Alexis Casas, Imperial College London

Co-Web Chair - Aaron Adler, BBN Technologies

Co-Web Chair - Prashant Vaidyanathan, Microsoft Research

Finance Chair - Traci Haddock-Angelli, iGEM Foundation

Virtual Chairs

Jet Mante, University of Colorado, Boulder

Dimitris Papamichail, The College of New Jersey

Andrea Cristina, Universidad de Ingeniería y Tecnología

Martin Gutierrez, Universidad Diego Portales

Lukas Buecherl, University of Colorado, Boulder

Alejandro Vignoni, Universitat Politècnica de València

Kenza Samlali, Concordia University

Samuel MD Oliveira, Boston University

Bio-Design Automation Consortium

President - Aaron Adler, BBN Technologies

Vice-President - Natasa Miskov-Zivanov, Carnegie Mellon University

Treasurer - Traci Haddock, iGEM Foundation

Board Member - Douglas Densmore, Boston University

Program Committee

Aaron Adler	BBN Technologies
Jacob Beal	BBN Technologies
Lukas Buecherl	CU Boulder
Alexis Casas	Imperial College London
Zak Costello	Generate Biomedicines
Martín Gutiérrez	Universidad Diego Portales
Shruti Joshi	MIT-ADT university
Ernst Oberortner	DOE Joint Genome Institute
Luis Ortiz	Boston University
Dimitris Papamichail	The College of New Jersey
Marilene Pavan	Lanzatech
William Poole	Caltech
Kenza Samlali	Concordia University
Radhakrishna Sanka	Boston University
Jenhan Tao	Generate Biomedicines
Prashant Vaidyanathan	Microsoft Research
Eric Young	Worcester Polytechnic Institute
Zhen Zhang	Utah State University

Program

All times BST. Conference to be conducted over Zoom.

Monday, 20th September 2021 15:00 - 15:15 Welcome & Opening Remarks Prashant Vaidyanathan

15:15 - 16:45 **Session 1: Standardization of Biological Components**, Chair: Andrea Cristina

- 15:15-15:30 *Network visualisation of synthetic biology designs*
Matthew Crowther, Anil Wipat and Ángel Goñi-Moreno
- 15:30-15:55 *Data Representation in the DARPA SD2 Program*
Nicholas Roehner, Jacob Beal, Bryan Bartley, Richard Markeloff, Tom Mitchell, Tramy Nguyen, Daniel Sumorok, Nicholas Walczak, Chris Myers, Zach Zundel, James Scholz, Benjamin Hatch, Mark Weston and John Colonna-Romano
- 15:55-16:20 *Excel-SBOL Converter: Creating SBOL from Excel Templates and Vice Versa*
Julian Abam, Jeanet Mante, Isabel Pöttsch, Jake Beal and Chris Myers
- 16:20-16:45 *Towards collaborative and automated development of resources for data standards in synthetic biology*
Jake Sumner Ajibode, Jacob Beal, James Scott-Brown, Thomas Gorochowski, Chris Myers and Goksel Misirli

16:45 - 17:00 **Short Break**

17:00 - 18:00 **Live Keynote I: Dr. Tijana Radivojevic**

- **Title:** Guiding synthetic biology via machine learning and multi-omics technologies.
- **Abstract:** Synthetic biology allows us to bioengineer cells to synthesize novel valuable molecules such as renewable biofuels or anticancer drugs. However, traditional synthetic biology approaches involve ad-hoc engineering practices, which lead to long development times. One of the most important challenges in bioengineering is effectively using multi-omics data to guide metabolic engineering towards higher production levels. In this talk, I will show our efforts in developing pipelines for collection of multi-omics datasets, their analysis through machine learning, and the production of recommendations, with the goal of accelerating the Design-Build-Test-Learn (DBTL) cycle. I will focus on machine learning techniques that, trained on the multi-omics datasets, provide actionable recommendations predicted to optimize strain performance and increase production through several DBTL cycles. These tools help guide synthetic biology in a systematic fashion, without the need for a full mechanistic understanding of the biological system. Our tools also aim to enable commercially-relevant bioengineering and are currently being deployed by industrial and academic partners.

18:00 - 18:30 **Gather Social Hour**

18:30 - 20:30 **Workshop 1: SBOL Version 3: Data Exchange throughout the Bioengineering Lifecycle**

Tuesday, 21st September 2021

15:00 - 15:05 **Welcome & Opening Remarks** Jenhan Tao

15:05 - 16:00 **Session 2: Screening Methods**, Chair: Kenza Samlali

- 15:05-15:30 *The Bioware Cyber-Fluidic Platform: A Holistic Approach to Digital Microfluidics*
Georgi Tanev, Luca Pezzarossa and Jan Madsen
- 15:30-15:45 *An Investigative Platform Comprising Cell-Free Transcription- Translation and Electron Microscopy for Studying Bacteriophages*
Joseph Wheatley, Sahan Liyanagedera, Ian Hands-Portman, Antonia Sagona and Vishwesh Kulkarni
- 15:45-16:00 *Engineering SpyTag Bacteriophage K1F for Directional Immobilisation*
Sahan Liyanagedera, Joseph Wheatley, Alyona Biketova, Ian Hands-Portman, Antonia Sagona, Kevin Purdy, Tamas Feher and Vishwesh Kulkarni

16:00 - 16:15 **Short Break**

16:15 - 17:15 **Live Keynote II: Dr. Thomas E. Goroehowski**

- **Title:** Programming biology – controlling the flow of molecular machines empowering life.
- **Abstract:** Synthetic genetic circuits are composed of many interconnected parts that must control the flows of transcriptional and translational machinery such that a desired biological computation can be implemented. A major challenge when developing such circuits is that the genetic parts used often display unexpected changes in their behavior when pieced together in new ways. Such changes can arise due to contextual effects or unintended interactions with the host cell. In this talk, I will demonstrate how we have been using a variety of sequencing technologies to create a genetic debugger to pinpoint the root of such failures, as well as our recent efforts to develop “tunable” genetic parts whose functions can be dynamically altered to fix many of these common issues. I will also discuss some of our recent efforts to consider the role of evolution in biological design and the concept of the ‘evotype’ as a way to reason about the evolutionary potential of engineered biology. Taken together, our work provides a more complete and quantitative view of the inner workings of genetic circuits, offers a route to engineering more robust and adaptive functionalities in living cells, and improves our understanding of the rules governing the effective reprogramming of biology.

17:15 - 17:45 **Gather Social Hour**

17:45 - 19:45 **Workshop 2: Part I -From Chemical Reaction Network Compilation to Bayesian Parameter Inference (Compilation with BioCRNpyler)**

19:45 - 21:45 **Workshop 2: Part II -From Chemical Reaction Network Compilation to Bayesian Parameter Inference (Model Reduction and Inference)**

Wednesday, 22nd September 2021

15:00 - 15:05 **Welcome & Opening Remarks** Marilene Pavan

15:05 - 15:55 **Session 3: Metabolic & Knowledge Engineering I**, Chair: Alejandro Vignoni

- 15:05 - 15:30 *Modeling of the engineered production of curcumin in Escherichia coli*

Michael Cotner, Ellie Brown, Jixun Zhan and Zhen Zhang

- 15:30 - 15:55 *Codon-Optimized Degenerate Codon Set Design Tool*

Akira Takada, Tomer Aberbach, Nicholas Carpino, Georgios Papamichail and Dimitris Papamichail

15:55 - 16:00 **Short Break**

16:00 - 17:30 **Workshop 3: Machine Learning Aided Advances in Synthetic Biology**

17:30 - 17:45 **Short Break**

17:45 - 18:50 **Session 4: Metabolic & Knowledge Engineering II**, Chair: Dimitris Papamichail

- 17:45 - 18:00 *iBioSim Server: a Tool for Improving the Workflow for Genetic Design and Modeling*

Thomas Stoughton, Lukas Buecherl, Payton Thomas, Pedro Fontanarrosa and Chris Myers

- 18:00 - 18:25 *Optimizing and Classifying Literature Events for Automated Model Extension*

Casey Hansen, Julia Kisslinger, Neal Krishna, Emilee Holtzapple, Yasmine Ahmed and Natasa Miskov-Zivanov

- 18:25 - 18:50 *New advances in the automation of context-aware information selection and guided model assembly*

Yasmine Ahmed, Adam A. Butchy, Khaled Sayed, Cheryl Telmer and Natasa Miskov-Zivanov

18:50 - 19:15 **Gather Social Hour**

19:15 - 21:15 **Workshop 4: DIY Microfluidics CAD - Extending 3DuF for fun and publications**

Thursday, 23rd September 2021

15:00 - 15:05 **Welcome & Opening Remarks** Alexis Casas

15:05 - 16:45 **Session 5: Computer-aided design, modelling, and simulation**, Chair: Martin G. Pescarmona

- 15:05 - 15:30 *A Comparison of Weighted Stochastic Simulation Methods*

Payton Thomas, Mohammad Ahmadi, Hao Zheng and Chris Myers

- 15:30 - 15:55 *LOICA: Logical Operators for Integrated Cell Algorithms*

Gonzalo Vidal, Guillermo Yáñez-Feliú, Carlos Vidal-Céspedes and Timothy James Rudge

- 15:55 - 16:20 *Biophysical Technology Mapping of Genetic Circuits*

Nicolai Engelmann, Tobias Schladt, Erik Kubaczka, Christian Hochberger and Heinz Koepl

- 16:20 - 16:45 *Automated translation of logical models to SystemVerilog enables simulation speedup*

Eric Li, Emilee Holtzapple, Niteesh Sundaram and Natasa Miskov-Zivanov

16:45 - 17:15 **Break / Short Discussion - Distributed mini-biofoundries**

17:15 - 21:15 **Workshop 5: Flapjack: Data Management and Analysis for Genetic Circuit Characterization**

Friday, 24th September 2021

14:00 - 14:05 **Welcome & Opening Remarks**

14:05 - 14:50 **Workshop 6: BMSS - An Automated BioModel Selection System for Gene Circuit Designs**

14:50 - 15:00 **Short Break**

15:00 - 15:40 **Session 6: Machine Learning**, Chair: Jeanet Mante

- 15:00 - 15:15 *Improving predictability in bio-design using ensemble models*
Bret Peterson, Tijana Radivojevic and Hector Garcia Martin
- 15:15 - 15:40 *Comparison of Extrinsic and Intrinsic Noise Model Predictions for Genetic Circuit Failures*
Pedro Fontanarrosa, Lukas Buecherl and Chris J. Myers

15:40 - 16:00 **GatherSocial Hour**

16:00 - 18:00 **Workshop 7: Visualizing biological designs using SBOL visual**

18:00 - 18:15 **Short Break**

18:15 - 18:45 **Session 7: Data Repositories**, Chair: Samuel MD Oliveira

- 18:15 - 18:30 *SynBioHub2 - Providing an Intuitive and Maintainable Genetic Design Repository*
Benjamin Hatch, Jeanet Mante, Chris Myers and Eric Yu
- 18:30 - 18:45 *A database for ligand-inducible genetic biosensors*
Simon d'Oelsnitz

18:45 - 20:45 **Workshop 8: Principles of genetic circuit design: programming living cells to perform novel functions**

20:45 - 21:00 **Closing Remarks**

Keynote Presentation

Guiding synthetic biology via machine learning and multi-omics technologies

Tijana Radivojevic



Speaker Biography

Dr. Tijana Radivojevic is a Data Scientist at Agile BioFoundry and Joint BioEnergy Institute, Lawrence Berkeley National Lab (LBNL). Her current interests lie in helping bioengineering become a mature engineering discipline. She has been working on development of machine learning based algorithms and tools for guiding and predicting outcomes of bioengineering, while capturing the associated uncertainty. Prior to joining LBNL, her research focused on development of methodologies in computational statistics, applied across domains such as finance, reservoir simulation, molecular simulation, breast cancer. Tijana holds a PhD degree in Applied Mathematics from the University of the Basque Country, Spain, and MSc, BSc degrees in Financial Mathematics from the University of Novi Sad, Serbia.

Keynote Abstract

Synthetic biology allows us to bioengineer cells to synthesize novel valuable molecules such as renewable biofuels or anticancer drugs. However, traditional synthetic biology approaches involve ad-hoc engineering practices, which lead to long development times. One of the most important challenges in bioengineering is effectively using multi-omics data to guide metabolic engineering towards higher production levels. In this talk, I will show our efforts in developing pipelines for collection of multi-omics datasets, their analysis through machine learning, and the production of recommendations, with the goal of accelerating the Design-Build-Test-Learn (DBTL) cycle. I will focus on machine learning techniques that, trained on the multi-omics datasets, provide actionable recommendations predicted to optimize strain performance and increase production through several DBTL cycles. These tools help guide synthetic biology in a systematic fashion, without the need for a full mechanistic understanding of the biological system. Our tools also aim to enable commercially-relevant bioengineering and are currently being deployed by industrial and academic partners.

Keynote Presentation

Programming biology – controlling the flow of molecular machines empowering life

Thomas E. Gorochofski



Speaker Biography

Dr. Thomas E. Gorochofski is a Royal Society University Research Fellow at the University of Bristol and Co-Director of the Bristol BioDesign Institute (BBI). His laboratory focuses on exploring the molecular and biophysical mechanisms that individual cells and groups of cells use to make sense of their world and process information. By applying tools from the field of synthetic biology to create new genetic systems from the ground-up, his laboratory then probes these artificial systems using novel techniques based on diverse content-rich sequencing methods and advanced computer models, with the aim of better understanding the rules governing how biological parts are best pieced together to perform useful computations. Elucidating the computational architecture of living cells and cellular collectives opens new ways of effectively reprogramming them to tackle problems spanning the sustainable production of materials to novel therapeutics, while also providing fundamental insight into how biology orchestrates the complex processes and structures sustaining life.

Keynote Abstract

Synthetic genetic circuits are composed of many interconnected parts that must control the flows of transcriptional and translational machinery such that a desired biological computation can be implemented. A major challenge when developing such circuits is that the genetic parts used often display unexpected changes in their behavior when pieced together in new ways. Such changes can arise due to contextual effects or unintended interactions with the host cell. In this talk, I will demonstrate how we have been using a variety of sequencing technologies to create a genetic debugger to pinpoint the root of such failures, as well as our recent efforts to develop “tunable” genetic parts whose functions can be dynamically altered to fix many of these common issues. I will also discuss some of our recent efforts to consider the role of evolution in biological design and the concept of the ‘evotype’ as a way to reason about the evolutionary potential of engineered biology. Taken together, our work provides a more complete and quantitative view of the inner workings of genetic circuits, offers a route to engineering more robust and adaptive functionalities in living cells, and improves our understanding of the rules governing the effective reprogramming of biology.

Regular Talks

1	Towards collaborative and automated development of resources for data standards in synthetic biology <i>Jake Sumner Ajibode, Jacob Beal, James Scott-Brown, Thomas Gorochowski, Chris Myers and Goksel Misirli</i>	14
2	A Comparison of Weighted Stochastic Simulation Methods <i>Payton Thomas, Mohammad Ahmadi, Hao Zheng and Chris Myers</i>	17
3	The Bioware Cyber-Fluidic Platform: A Holistic Approach to Digital Microfluidics <i>Georgi Tanev, Luca Pezzarossa and Jan Madsen</i>	20
4	Comparison of Extrinsic and Intrinsic Noise Model Predictions for Genetic Circuit Failures <i>Pedro Fontanarrosa, Lukas Büecherl and Chris J. Myers</i>	25
5	Modeling of the engineered production of curcumin in <i>Escherichia coli</i> <i>Michael Cotner, Ellie Brown, Jixun Zhan and Zhen Zhang</i>	29
6	Biophysical Technology Mapping of Genetic Circuits <i>Nicolai Engelmann, Tobias Schladt, Erik Kubaczka, Christian Hochberger and Heinz Koeppel</i>	33
7	Automated translation of logical models to SystemVerilog enables simulation speedup <i>Eric Li, Emilee Holtzapple, Nitesh Sundaram and Natasa Miskov-Zivanov</i>	36
8	Excel-SBOL Converter: Creating SBOL from ExcelTemplates and Vice Versa <i>Julian Abam, Jeanet Mante, Isabel Pöttsch, Jake Beal and Chris Myers</i>	40
9	LOICA: Logical Operators for Integrated Cell Algorithms <i>Gonzalo Vidal, Guillermo Yáñez-Feliú, Carlos Vidal-Céspedes and Timothy James Rudge</i>	44
10	Optimizing and Classifying Literature Events for Automated Model Extension <i>Casey Hansen, Julia Kisslinger, Neal Krishna, Emilee Holtzapple, Yasmine Ahmed and Natasa Miskov-Zivano</i>	49
11	New advances in the automation of context-aware information selection and guided model assembly <i>Yasmine Ahmed, Adam A. Butchy, Khaled Sayed, Cheryl Telmer and Natasa Miskov-Zivanov</i>	52
12	Codon-Optimized Degenerate Codon Set Design Tool <i>Akira Takada, Tomer Aberbach, Nicholas Carpino, Georgios Papamichail and Dimitris Papamichail</i>	56
13	Data Representation in the DARPA SD2 Program <i>Nicholas Roehner, Jacob Beal, Bryan Bartley, Richard Markeloff, Tom Mitchell, Tramy Nguyen, Daniel Sumorok, Nicholas Walczak, Chris Myers, Zach Zundel, James Scholz, Benjamin Hatch, Mark Weston and John Colonna-Romano</i>	59

Short Talks

1	Network visualisation of synthetic biology designs <i>Matthew Crowther, Anil Wipat and Ángel Goñi-Moreno</i>	63
2	A database for ligand-inducible genetic biosensors <i>Simon d'Oelsnitz</i>	66
3	Improving ensemble model predictions in a general biosynthesis experiment modeling tool <i>Bret Peterson, Tijana Radivojevic and Hector Garcia Martin</i>	68
4	iBioSim Server: a Tool for Improving the Workflow for Genetic Design and Modeling <i>Thomas Stoughton, Lukas Buecherl, Payton Thomas, Pedro Fontanarroza and Chris Myers</i>	72
5	SynBioHub2 - Providing an Intuitive and Maintainable Genetic Design Repository <i>Benjamin Hatch, Jeanet Mante, Chris Myers and Eric Yu</i>	75
6	An Investigative Platform Comprising Cell-Free Transcription- Translation and Electron Microscopy for Studying Bacteriophages <i>Joseph Wheatley, Sahan Liyanagedera, Ian Hands-Portman, Antonia Sagona and Vishwesh Kulkarni</i>	79
7	Engineering SpyTag Bacteriophage K1F for Directional Immobilisation <i>Sahan Liyanagedera, Joseph Wheatley, Alona. Yu Biketova, Ian Hands-Portman, Antonia Sagona, Kevin Purdy, Tamas Feher and Vishwesh Kulkarni</i>	81

Towards collaborative and automated development of resources for data standards in synthetic biology

Jake Sumner Ajibode¹, Jacob Beal², James Scott-Brown³ Thomas E. Gorochowski⁴, Chris J. Myers⁵, Göksel Mısırlı¹

¹Keele University, ²Raytheon BBN Technologies, ³University of Oxford, ⁴University of Bristol, ⁵University of Colorado Boulder

w8q52@students.keele.ac.uk, jakebeal@ieee.org, james@jamesscottbrown.com
thomas.gorochowski@bristol.ac.uk, chris.myers@colorado.edu, g.misirli@keele.ac.uk

1 INTRODUCTION

Data standards in synthetic biology are becoming increasingly important as the number of tools addressing different needs grows. This covers the design of genetic circuits and visualizing and storing underlying information regarding biological designs. The Synthetic Biology Open Language (SBOL) [5] has been developed to provide a mechanism to capture and exchange a shared understanding of such designs and related information. Moreover, SBOL Visual [1] has been developed to standardize the representation of biological designs via well-defined glyphs and rules for how these are combined and connected.

The SBOL community has adopted GitHub for the collaborative development of specifications online. The community website is built using GitHub pages and is accessible via HTTP. The SBOL data and visual standards follow an incremental release process by applying common Git-related processes such as branching for incremental developments and making releases with explicit version numbering. Standard enhancement proposals are prepared by community members, discussed, and finally incorporated into minor or major SBOL versions. The changes between versions can affect existing tools, and additional mechanisms are needed for documentation and verification. As the number and frequency of these new releases increase, there is a growing desire to automate the preparation of documentation and other resources produced. Such an automated process would also open opportunities to allow for the testing of changes before incorporating them into formal releases.

An automated approach is of particular interest to the SBOL Visual standard that has grown rapidly over the past few years to contain a large number of glyphs and associated data. Each glyph definition comes with a Scalable Vector Graphics (SVG) file, providing information about how the glyph can be connected to the others in a biological design. Complementary metadata about each glyph is stored in a markdown file to ensure this information is human-readable. Each markdown file includes information about the types of a biological part or interaction that the glyph can be used with via ontological terms from the Sequence Ontology (SO)

[4] and the Systems Biology Ontology (SBO) [3]. In addition, these files include lists of allowed alternative glyphs, free-text notes and other related information. These markdown files should ideally be validated after each change in the context of all other glyph files and metadata from different files to prevent human error.

Previously, we developed SBOL Visual Ontology version 2 (SBOL-VO2) [7] to make the metadata about these glyphs machine accessible. We also developed the visual ontology web service (VOVS) to make the glyphs searchable and to find the most relevant glyph for a provided SO or SBO term representing a biological part or interaction. An ontological version of SBOL 2 has also been developed [6]. However, the web service is currently not suitable for supporting the automation of new releases or for testing it against a specific snapshot of the standard. Moreover, these resources need to be updated to incorporate the new SBOL 3 data and visual standards.

Here, we report the latest developments related to the automation of the SBOL Visual standard. This work involves developing a new SBOL Visual ontology for version 3 (SBOL-VO3) and improving the visual ontology web service for a better user experience. We then demonstrate an automated workflow to generate the visual ontology and related documentation. This will accelerate the ability to release a new SBOL Visual specification and enable error checking to improve the robustness of the process.

2 SBOL VISUAL ONTOLOGY 3 (SBOL-VO3)

In the new version of the SBOL Visual ontology, ontological terms related to glyphs have been grouped into four main branches to represent sequence features, molecular species, molecular interactions for binary relationships, and biological processes with multiple inputs and outputs. Sequence feature terms are linked to respective SO terms, while SBO is used in all other terms to indicate the types of all molecular interactions, processes and non-DNA-based molecular terms. Another new feature is to capture a network-level representation of biological designs using directed edges. SBO terms are used to provide information about the roles of molecules

in molecular interactions and processes, and hence the direction of edges. The `hasHead` and `hasTail` terms indicate the direction of binary interactions. For example, ‘`hasHead some (role some SBO:0000642)`’ axiom for an inhibition interaction specifies that the direction of the interaction starts from an entity with the inhibitor role (SBO:0000642). Similarly, `hasIncoming` and `hasOutgoing` properties indicate the inputs and outputs for participating biological molecules. This visual ontology refers to terms from the latest SBOL 3 data ontology (SBOL-OWL3).

In line with the major changes in SBOL 3, we developed the proof-of-concept SBOL-OWL3 ontology. This ontology captures the rich SBOL 3 data model, which uses a graph-based approach to represent biological designs. Hence, SBOL-VO3 and SBOL-OWL3 ontologies act as machine-accessible sources to connect the SBOL Visual glyphs with SBOL specific data entities.

3 SBOL VISUAL ONTOLOGY WEB SERVICE

We added new features to the previously developed SBOL Visual Ontology Web Service [7]. Specifically, we decoupled the web service from a particular version of an SBOL Visual ontology and other related SO and SBO ontologies. These external resources are now loaded directly from the web. Different versions of the visual ontology are held at a GitHub repository. By default, the web service uses the latest version of the SBOL Visual ontology and the latest glyph files. However, users can also specify a particular release number. Alternatively, users can provide HTTP URLs for the ontology and the base glyph folder. These URLs can also be for a specific GitHub branch to test resources under development, and different media types for glyphs can be requested.

4 AUTOMATION

Our goal is to develop a fully automated workflow to test incremental or major updates and use them directly via the SBOL Visual web service (Figure 1). To facilitate this approach, we aimed to use GitHub Actions. As a result, specific GitHub-related events such as new releases or commits to a particular branch can trigger our workflow to auto-generate the SBOL Visual ontology, related documentation and make the visual glyphs directly available via the web service. We developed scripts that can be linked to GitHub actions so that the SBOL community can take advantage of these resources quickly.

5 CONCLUSION

Ontologies are ideal for capturing domain knowledge. We used this approach to provide a machine-accessible representation of the SBOL 3 data and visual standards. In order to help the SBOL community deal with incremental and major changes, we adopt an automated approach based on GitHub

actions. This automation enables testing changes rapidly and releasing them to the community quickly. The improvements enhance the user experience via different additional options when using the web service. The SBOL community is working on a standard format for parametric glyphs that can be customized for different properties such as color, width, and height [2]. Our goal is to incorporate additional metadata about these parameters into the SBOL Visual ontology and use the web service to serve parametric glyphs for design and layout tools in the future. SBOL Visual glyphs, the automation workflow and the SBOL Visual Ontology are available from <https://github.com/SynBioDex/SBOL-visual> and <https://github.com/SynBioDex/sbol-visual-ontology>.

6 ACKNOWLEDGEMENTS

We thank SBOL Industrial Consortium for funding this work and the SBOL community for support and feedback. This document does not contain technology or technical data controlled under either U.S. International Traffic in Arms Regulation or U.S. Export Administration Regulations.

REFERENCES

- [1] BAIG, H., FONTANAROSSA, P., KULKARNI, V., McLAUGHLIN, J., VAIDYANATHAN, P., BARTLEY, B., BHAKTA, S., BHATIA, S., BISSELL, M., CLANCY, K., COX, R. S., MORENO, A. G., GOROCHOWSKI, T., GRUNBERG, R., LEE, J., LUNA, A., MADSEN, C., MISIRLI, G., NGUYEN, T., NOVERE, N. L., PALCHICK, Z., POCOCK, M., ROEHNER, N., SAURO, H., SCOTT-BROWN, J., SEXTON, J. T., STAN, G.-B., TABOR, J. J., TERRY, L., VILAR, M. V., VOIGT, C. A., WIPAT, A., ZONG, D., ZUNDEL, Z., BEAL, J., AND MYERS, C. Synthetic biology open language visual (sbol visual) version 2.3. *J. Integr. Bioinform.* (2021), 20200045.
- [2] CLARK, C. C., SCOTT-BROWN, J., AND GOROCHOWSKI, T. E. *parasbol*: a foundation for standard-compliant genetic design visualization tools. *Synthetic Biology* (2021).
- [3] COURTOT, M., JUTY, N., KNÜPFER, C., WALTEMATH, D., ZHUKOVA, A., DRÄGER, A., DUMONTIER, M., FINNEY, A., GOLEBIEWSKI, M., HASTINGS, J., HOOPS, S., KEATING, S., KELL, D. B., KERRIEN, S., LAWSON, J., LISTER, A., LU, J., MACHNE, R., MENDES, P., POCOCK, M., RODRIGUEZ, N., VILLEGER, A., WILKINSON, D. J., WIMALARATNE, S., LAIBE, C., HUCKA, M., AND LE NOVÈRE, N. Controlled vocabularies and semantics in systems biology. *Mol. Syst. Biol.* 7, 1 (2011), 543.
- [4] EILBECK, K., LEWIS, S. E., MUNGALL, C. J., YANDELL, M., STEIN, L., DURBIN, R., AND ASHBURNER, M. The sequence ontology: a tool for the unification of genome annotations. *Genome Biol.* 6, 5 (2005), R44.
- [5] McLAUGHLIN, J. A., BEAL, J., MISIRLI, G., GRÜNBERG, R., BARTLEY, B. A., SCOTT-BROWN, J., VAIDYANATHAN, P., FONTANAROSSA, P., OBERORTNER, E., WIPAT, A., GOROCHOWSKI, T. E., AND MYERS, C. J. The synthetic biology open language (sbol) version 3: Simplified data exchange for bioengineering. *Front Bioeng. Biotechnol.* 8 (2020), 1009.
- [6] MISIRLI, G., TAYLOR, R., GONI-MORENO, A., McLAUGHLIN, J. A., MYERS, C., GENNARI, J. H., LORD, P., AND WIPAT, A. *Sbol-owl*: An ontological approach for formal and semantic representation of synthetic biology information. *ACS Synth. Biol.* 8, 7 (2019), 1498–1514.
- [7] MISIRLI, G., BEAL, J., GOROCHOWSKI, T. E., STAN, G.-B., WIPAT, A., AND MYERS, C. J. *Sbol visual 2 ontology*. *ACS Synth. Biol.* 9, 4 (04 2020), 972–977.

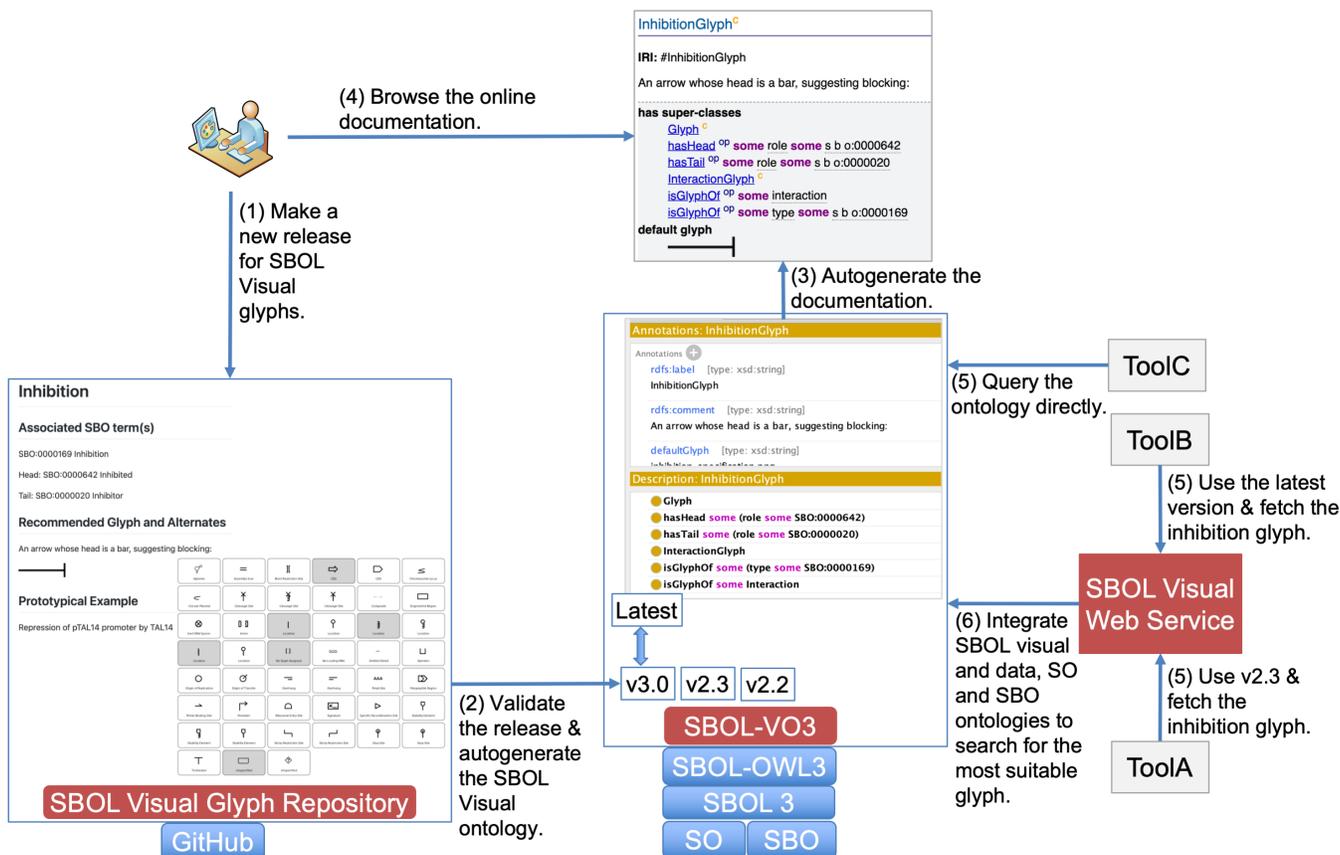


Figure 1: A use case involving the use of the automated approach. A user modifies the metadata for an SBOL Visual glyph and makes a new release. This release triggers the automatic regeneration of the SBOL Visual ontology and the online documentation. Tools can start using this change immediately. Alternatively, they can continue using a specific or an older version of the ontology and glyphs. Tools can either query the ontology file or utilize the web service to find the most suitable glyphs for given SO or SBO terms.

A Comparison of Weighted Stochastic Simulation Methods

Payton J Thomas¹, Mohammad Ahmadi², Hao Zheng², Chris J. Myers³,
¹University of Utah, ²University of South Florida, ³University of Colorado Boulder
 chris.myers@colorado.edu

1 INTRODUCTION

Despite occurring with low frequency, rare events can have devastating effects on biological systems. For example, rare biochemical events have been demonstrated to contribute to cancerous phenotypes by inactivating tumor-suppressing genes [1]. It is therefore important that computational methods be developed to analyze the probability of rare events.

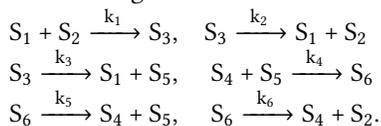
Exact trajectories of biochemical reaction networks may be determined with *molecular dynamics*, wherein, given the initial position and momentum of each atom in the system, the complete state of the system can be determined at any time [4]. Unfortunately, such methods are computationally intractable for most systems. Instead, *stochastic chemical kinetics* (SCK) may be used to generate many potential trajectories for a system and approximate the probability of some event occurring [8].

Rare events can be problematic for stochastic simulation because the number of trajectories that must be generated to approximate the probability of a rare event may be computationally prohibitive. To address this issue, a variety of stochastic simulation algorithms have been developed that utilize *importance sampling* (IS) techniques to better estimate the probability of rare events [3, 6, 7]. In this abstract, three such algorithms are examined to determine how well they address the problem of rare event simulation.

The first algorithm that will be examined is the *weighted stochastic simulation algorithm* (wSSA) [6], which first applied IS techniques to biochemical network simulation. The second algorithm that will be examined is the *state-dependent biasing method for importance sampling* (swSSA) [7]. The third algorithm that will be examined is the *guided weighted stochastic simulation algorithm* (guided wSSA) [3].

2 RESULTS

The efficacy of each stochastic simulation method was tested on a six-reaction model of a biochemical futile cycle. This network is given as follows:



where

$$k_1 = k_2 = k_4 = k_5 = 1, \quad k_3 = k_6 = 0.1.$$

In the model, the initial state is

$$\begin{aligned} X_1(0) = X_4(0) = 1, \quad X_2(0) = X_5(0) = 50, \\ X_3(0) = X_6(0) = 0. \end{aligned}$$

The rare event of interest is $X_5 \rightarrow 40$ within 100 time units, which is unlikely because the symmetry of the initial molecule counts and reaction rate constants will keep the system near the initial state with high probability. Futile cycles of this kind exist biologically in GTPase cycles, MAPK cascades, and glucose mobilization [2].

Rare events are difficult to simulate because the number of traditional SSA runs necessary to see a rare event of interest occur even once can be very high. Kuwahara and Mura solve this issue by increasing the likelihood of certain reactions occurring in simulation and decreasing the likelihood of others. Each run is then assigned a weight specific to the sequence of reactions that occurred such that the mean run weight is a sample estimator for the probability of the rare event of interest. This method requires that the user manually input the IS biasing parameters that are applied to each reaction.

In the small six-reaction example network, reaction three produces species five and reaction six consumes species five, so reaction three must be biased downward and reaction six must be biased upward. To this end, a single biasing parameter $0 < \delta$ was introduced such that the rate of reaction three is multiplied by δ and the rate of reaction six is divided by δ . The performance of various magnitudes for δ is determined by comparing the true probability of a rare event $X_5 \rightarrow 40$ to the wSSA estimate after 10^2 runs for $0 < \delta \leq 1.5$ with increment 0.025 (Figure 1(a)).

As the species population changes throughout the course of simulation, relative propensity of each reaction changes too. The key insight presented in [7] is that forcing a fixed IS biasing factor in wSSA will result in a narrow range of values that factor can take to produce an accurate estimate. That is because the fixed biasing factor must adjust relative propensities appropriately for most of the possible values they take throughout the simulation. Also, a fixed biasing parameter will increase/decrease relative propensity of a reaction at a state where it already has a high/low probability of selection, resulting in lower accuracy. Therefore, [7] introduces a biasing factor which is a function of the relative propensity of the reaction it is adjusting at the current state.

These functions are characterized by two sets of user inputs: (1) maximum amount of change allowed for each reaction, γ_j and (2) a threshold from which encouraging/discouraging reaction selection is stopped, ρ_{0j} .

Figure 2 shows the results of estimating the probability of the rare event of interest on six reaction network. Again, $R3$ is set to be biased downward and $R6$ is set to be biased upward. Fixing ρ_0 to be 0.6 for $R6$ and 0.2 for $R3$, and setting $\gamma_3 = \gamma_6 = \gamma_{max}$, true probability of the rare event is compared to swSSA estimates with $1 \leq \gamma_{max} \leq 5$ with increment 0.025 after 100 runs (Figure 1(b)).

To avoid reliance on user input and *a priori* knowledge of the system, Gillespie and Golightly calculate the conditioned expectation of reaction count over the remainder of the simulation for each reaction given that the rare state of interest is attained at the end of the simulation by assuming a constant reaction hazard and use that expectation to estimate an ideal amount of IS biasing [3, 5]. Unfortunately, the Guided wSSA may calculate a negative ideal biasing, and the resultant negative reaction rates cause errors in simulation. Inspection of the R code for the three example cases in Gillespie and Golightly reveals that a different method of dealing with these negatives is used in each case.

The Guided wSSA was ran using each of the three negative resolution methods with 10^3 runs to compare the performance of each method (Figure 1(c)).

3 DISCUSSION

The wSSA requires the user to select which reactions should be encouraged and which reactions should be discouraged. Although such a task might seem trivial for very simple models, deep insight into underlying dynamics of the network is necessary for more complex models. Also, the biasing factor for those selected reactions must be specified prior to simulation. This is a tricky task, since these parameters can arbitrarily take any value greater than zero and it is by no means obvious what values will result in accurate estimates just by considering the model. Moreover, the accuracy of the estimate is highly sensitive to these values. Selecting non-optimal biasing factors can result in an estimate even less accurate than one produced by running the original SSA for the same number of simulations.

The swSSA suffers from the same issues. Reactions which are to be encouraged/discouraged should be specified by the user. Furthermore, for each of those reactions, the maximum amount of change allowed as well as a threshold from which encouragement/discouragement should be applied must be set prior to simulation, resulting in twice as many parameters as the wSSA. Like with the wSSA, the accuracy of this method is sensitive to these parameters, although the swSSA generally produces more accurate estimates and shows more robustness against a wider range of these parameters.

The Guided wSSA eliminates the need for specifying a set of reactions to bias and parameter(s) associated with each of those reactions as in the wSSA and swSSA. Since, in the Guided wSSA, matrices are inverted to automatically recognize a suitable biasing factor, this method is inherently slower than the SSA, wSSA, and swSSA in simulating trajectories. This additional computational effort may be justified if the run weights have a variance which is considerably smaller than those produced by other methods (as is the case with experiments discussed in [3]). Estimation of the probability of the rare event discussed in Section 2 on the six reaction network model using guided wSSA produces a far less accurate estimate than estimating that with wSSA while setting $\delta = 0.6$. Running 200 simulations, it took guided wSSA 1.6 seconds to produce an estimate with the variance of 0.05 where it took wSSA 0.3 seconds to produce an estimate with the variance of 0.0015. The issue of complexity is demonstrated when the total runtimes of 10^5 runs of each algorithm are compared (Figure 1(d)).

In summary, the original wSSA may achieve rapid convergence and lower variance than competing methods, but only with a narrow set of biasing parameters that cannot be reliably determined for an arbitrary system. The swSSA demonstrates broader robustness to biasing variation, but estimates with a high proportional error with few runs, lessening its advantage over the SSA. The guided wSSA solves the issue of biasing parameter determination, but has poor run-time performance and converges slower than the wSSA with optimal biasing.

Acknowledgements

The authors of this work are supported by National Science Foundation Grant Nos. 1856740 and 1900542. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the funding agencies.

REFERENCES

- [1] ESTELLER, M. Epigenetics in cancer. *New England Journal of Medicine* 358, 11 (2008), 1148–1159.
- [2] FLOMENBOM, O., VELONIA, K., LOOS, D., MASUO, S., COTLET, M., ENGELBORGH, Y., HOFKENS, J., ROWAN, A. E., NOLTE, R. J., VAN DER AUWERAER, M., ET AL. Stretched exponential decay and correlations in the catalytic activity of fluctuating single lipase molecules. *Proceedings of the National Academy of Sciences* 102, 7 (2005), 2368–2372.
- [3] GILLESPIE, C. S., AND GOLIGHTLY, A. Guided proposals for efficient weighted stochastic simulation. *The Journal of chemical physics* 150, 22 (2019), 224103.
- [4] GILLESPIE, D. Handbook of materials modeling, chapter 5.11, 2005.
- [5] GOLIGHTLY, A., AND WILKINSON, D. J. Bayesian inference for markov jump processes with informative observations. *Statistical applications in genetics and molecular biology* 14, 2 (2015), 169–188.
- [6] KUWAHARA, H., AND MURA, I. An efficient and exact stochastic simulation method to analyze rare events in biochemical systems. *The Journal*

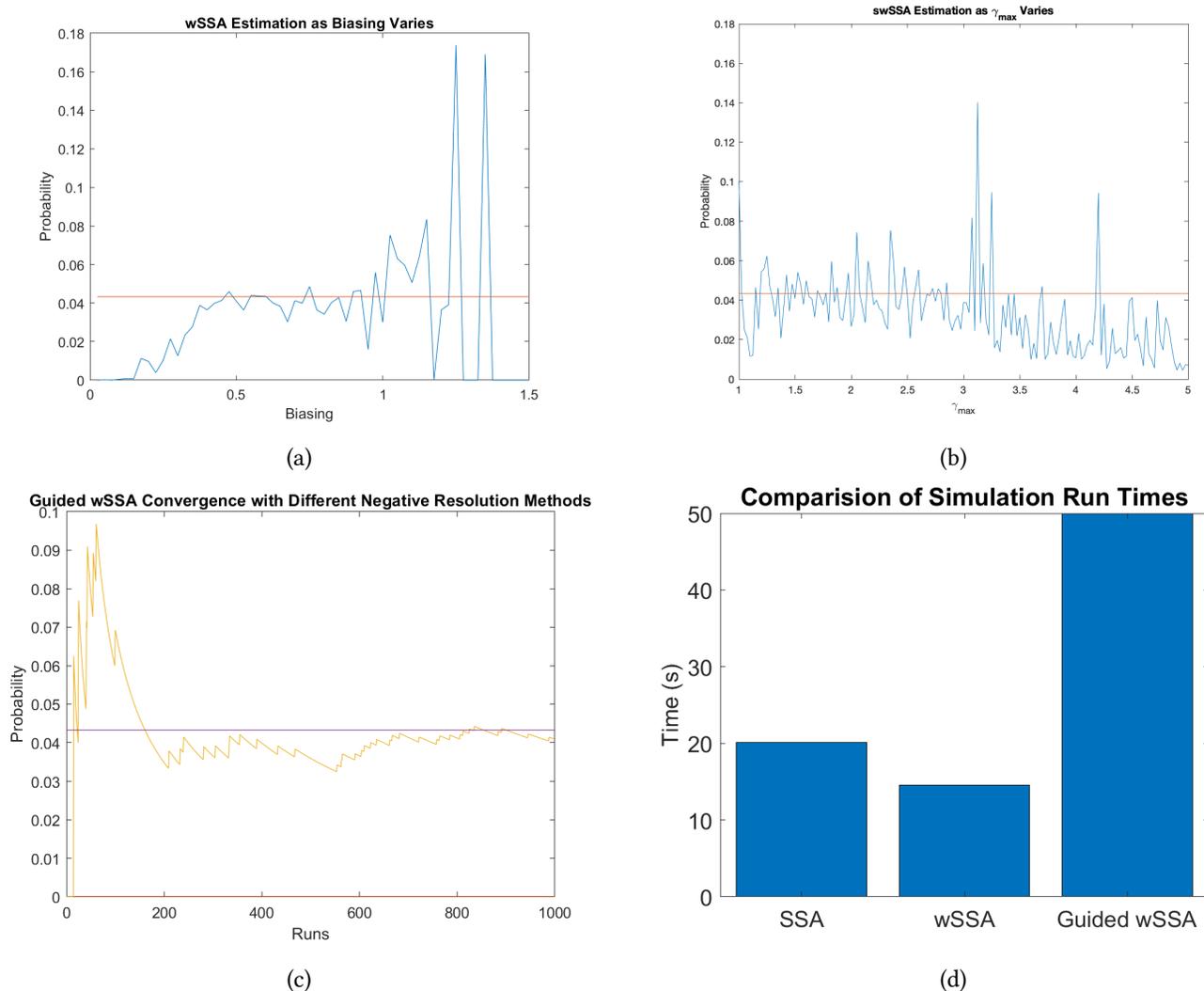


Figure 1: (a) True value of $P_{t \leq 100}(X_5 \rightarrow 40|x_0)$ (red) compared with the wSSA estimate at values of δ varying from 0.025 to 1.5 (blue). Note that $\delta = 1$ corresponds to the traditional SSA, and $\delta > 1$ corresponds to reciprocal weighting (decreases likelihood of reaching state of interest). (b) True value of $P_{t \leq 100}(X_5 \rightarrow 40|x_0)$ (red) compared with the swSSA estimate at values of γ_{max} varying from 1 to 5 (blue). (c) True value of $P_{t \leq 100}(X_5 \rightarrow 40|x_0)$ (purple) compared with the Guided wSSA estimate using each negative resolution method. Method A (yellow) and method C (blue) perform so similarly that method C is not visible. Method B fails to resolve negatives in general, and does not complete any runs. (d) The time to completion of 10^5 runs of each algorithm is compared. The relative computational complexity of the Guided wSSA makes it much slower than other methods. The wSSA (with an ideal biasing parameter) performs faster than the SSA despite performing more calculations because a large proportion of runs reach the state of interest before the total simulation time is reached.

of chemical physics 129, 16 (2008), 10B619.

- [7] ROH, M. K., GILLESPIE, D. T., AND PETZOLD, L. R. State-dependent biasing method for importance sampling in the weighted stochastic simulation algorithm. *The Journal of chemical physics* 133, 17 (2010), 174106.
 [8] SAMOILOV, M. S., AND ARKIN, A. P. Deviant effects in molecular reaction pathways. *Nature biotechnology* 24, 10 (2006), 1235–1240.

The Bioware Cyber-Fluidic Platform: A Holistic Approach to Digital Microfluidics

Georgi Tanev*, Luca Pezzarossa*, and Jan Madsen

Department of Applied Mathematics and Computer Science
Technical University of Denmark
Kongens Lyngby, Denmark
{geta,lpez,jama}@dtu.dk

1 INTRODUCTION

Digital microfluidic biochips have been in the research spotlight for over two decades bringing scientists on the quest of creating a universal portable bio-lab. Despite the unique first-order digital-to-fluidic control and numerous proof of concept applications, the digital biochips struggle to fulfil the promises of delivering extensive miniaturization, automation, and integration of processes used in the biomedical field. One of the main reasons is that although rather simple in construction, digital biochips require advanced hardware and software instrumentation, which poses many scientific, engineering, and integration challenges. Therefore, a holistic system approach is required to bridge the gap between the incremental technological developments and the vision for a digital microfluidic-based portable bio-lab. Inspired by the modern software-hardware symbiotic coexistence, we present the BiowareCFP – a cyber-fluidic platform aiming for seamless integration between the cyber, fluidic, and biological domains.

2 THE BIOWARECFP

The BiowareCFP is the natural continuation of our modular and reconfigurable digital microfluidic platform, which was presented at IWBD2019 [9]. The holistic approach to modularity allowed for extending the digital-biochip instrumentation along with the evolving requirement for biological sample handling. This was recently demonstrated by implementing the fundamental steps of two traditional, but rather different bioassays, namely full cell cloning, and magnetic beads-based ELISA [10]. Green fluorescent protein (GFP) was cloned into *E. coli* in a three-step process consisting of assembling recombinant DNA by PCR amplification, USER vector cloning, and transformation. These steps require temperature control, for which the heating capability of the digital biochip was utilized [9, 10]. Magnetic beads ELISA was demonstrated with both MRSA and SARS-CoV2 proteins, where the protocol required a magnetic field to capture and retain the magnetic beads for the washing steps [10].

In the early development stage of the BiowareCFP, the system was tested mainly with moving colored water droplets. Nevertheless, implementing the mentioned assays required handling of various fluids, including reagents with high enzymatic or protein content. Consequently, unspecific surface absorption became apparent in the form of skewing reagent concentrations or leading to reduced droplet mobility and consequential adhesion to the surface of the digital biochip. Moreover, keeping reagents at elevated temperatures for prolonged intervals inevitably leads to reagents evaporation and consequent bubble formation. Although techniques for mitigating these effects exist, their complete elimination is unlikely. Therefore, we are currently developing an extensive software control framework that implements real-time monitoring and dynamic experiment flow control for reliable handling of biological fluids.

Using digital biochips that can be reconfigured to different laboratory protocols requires a structured method to capture the protocol, verify the functionality through simulation, generate the platform commands, and implement real-time monitoring of the protocol execution. This process calls for a toolchain similar to the one used in classical software and hardware development, namely high- and low-level programming languages, a compiler, a debugger or simulator, and an execution environment. Although the majority of these components have already been researched separately, integration among them virtually does not exist resulting in digital biochips and instrumentation platforms that usually exhibit limited programming capabilities.

A comprehensive review of domain-specific languages for programmable biochemistry can be found in [7], while a variety of compilation techniques for digital microfluidics is outlined in [6]. Different approaches have been used in the development of these languages. For instance, BioCoder [3] extends the existing C++ language with a library enabling biologists to express the steps of a lab protocol. AquaCore [2] proposes an instruction set offering means to implement complete lab protocols on digital biochips, including operations such as merge, split, heat, etc. BioScript [5] is a standalone language characterized by an intuitive syntax optimized for human readability, as well as a type system ensuring that

*Joint main authorship. This research was funded by Novo Nordisk Fonden.

unsafe compounds interaction does not happen. State-of-the-art execution platforms for these domain specific-languages include DropBot [4], OpenDrop [1], Puddle [11], and the modular reconfigurable platform presented in [8]. Nevertheless, neither of the mentioned approaches fully account for dynamic protocol flow control based on real-time feedback from the digital biochip.

3 PROGRAMMING ENVIRONMENT

The proposed software framework aims to provide a set of tools that allows a user with little or no programming experience to design, compile, simulate, and execute bio protocols on the BiowareCFP. Particular emphasis is given to support real-time flow control on software and hardware levels which allow for decision making depending on sensing results and on-the-fly verification of the protocol execution. For example, in the case where a step of a protocol does not produce the expected result. In that case, intermediate actions can be taken to repeat or try to recover the step before proceeding, enabling real-time monitoring and correcting protocol execution. Figure 1 shows an overview of the main steps of the programming environment from the definition of the protocol to the execution on the digital biochip.

When programming for a biochip, the first step is to define the desired protocol using a high-level intuitive text or block-based graphical representation, which is then translated into the control flow graph of the protocol to be executed. For this, we are developing a web-based interface that allows to combine classic programming operations such as arithmetics, flow control, and printing, together with fluidic operations such as dispensing, mixing, incubation, temperature control, and sensing, as shown in Figure 3.

An abstract syntax tree of the protocol is then produced from the high-level representation and compiled into a target-specific representation. Information regarding the cyber-fluidic architecture, such as electrode and chip topology, availability and placement of actuators, sensors, etc. is provided as input to the compiler. The target representation of the compiled code is in BioAssembly - our domain-specific instruction set architecture (ISA) and assembly implementation. BioAssembly is inspired by classic computer ISAs, and it offers a platform-independent set of core instructions dedicated to digital biochips, as well as classic arithmetic, flow control, and memory access instructions from which complex droplet operations can be built. In addition, BioAssembly natively supports the execution of parallel synchronized tasks streamlining the droplet controls and interactions with the physical environment.

The compiled BioAssembly protocol is then executed by a virtual machine that emulates a processor-based system and interacts with the cyber-fluidic platform as a remote object through remote system calls. The BioAssembly instructions

that do not produce an action in the platform (e.g., arithmetic and logic, memory access, branches, etc.) are resolved and executed locally in the virtual machine, while the instructions that interact with the cyber-fluidic platform trigger the execution of the appropriate system calls. The use of a virtual machine decouples the BioAssembly execution and the functionality offered by the underlying cyber-fluidic platform since the low-level interaction with the physical platform needs to be resolved only once in the implementation of the action methods. The interface of the virtual machine is shown in Figure 4.

In addition to the execution on the BiowareCFP, the virtual machine can be connected to a simulator. The simulator is used to verify and debug the compiled protocol before running on the real platform. The simulator uses models of the cyber-fluidic platform actuators and sensors, as well as models of the droplets to simulate the protocol execution. In the future, we aim to also include unpredictable behavior using statistical models and handle in-droplet biochemical simulation.

4 VALIDATION AND EXPERIMENTS

The cyber-fluidic platform and its integration with the software framework were tested and validated by means of both artificial tests and traditional biological protocols.

Artificial tests were designed to validate specific features of the system. An example of this is a test where a continuous stream of droplets was sorted based on their color. Thus, demonstrating the execution of the code by the virtual machine and its capability of real-time flow control based on the color information captured by a color sensor integrated with the digital biochip. Figure 4 shows the interface of the virtual machine running the color sorting program and controlling its execution on the physical cyber-fluidic platform.

The tested real-life protocols include magnetic beads-based ELISA (for MRSA and SARS-CoV2 proteins) and the fundamental steps used in full cell cloning, where details can be found in [10]. Figure 2 shows an instance of the BiowareCFP configured for temperature calibration of the three individual temperature zones for PCR. After calibration, a space domain PCR was performed by moving and keeping the PCR mix droplet for predefined times at the annealing, melting, and elongation zones. Currently, these tests are not fully automated, but require supervision and sporadic interventions to mitigate bubble formation and bio-fouling. Nevertheless, they prove that real-life protocols can be entirely executed on our platform.

Future work includes further development of the software framework user interface, sensor and actuator integration on the digital biochip, and further testing of traditional protocols on the BiowareCFP such as sample preparation for single-cell proteomics.

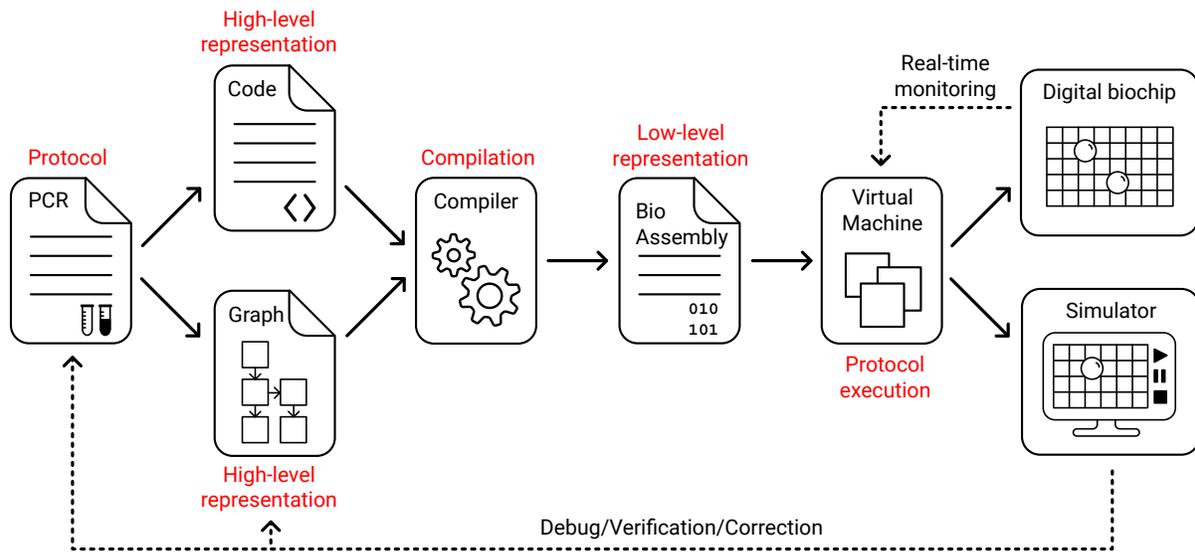


Figure 1: Overview of the main steps of the programming environment starting from the definition of the protocol and terminating with the execution on the biochip. The simulator is used to debug, verify, and correct the protocol.

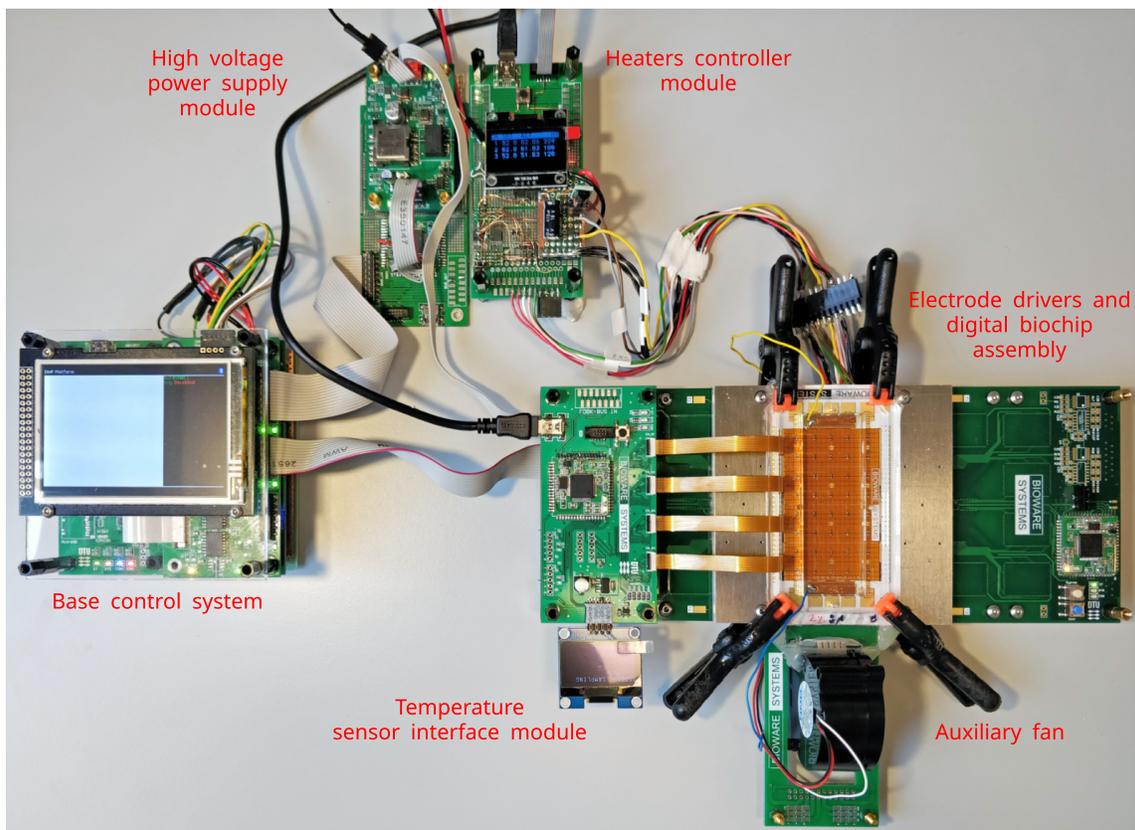


Figure 2: The BiowareCFP configured for temperature calibration of three temperature regions used for the annealing, melting, and elongation steps of the PCR protocol.

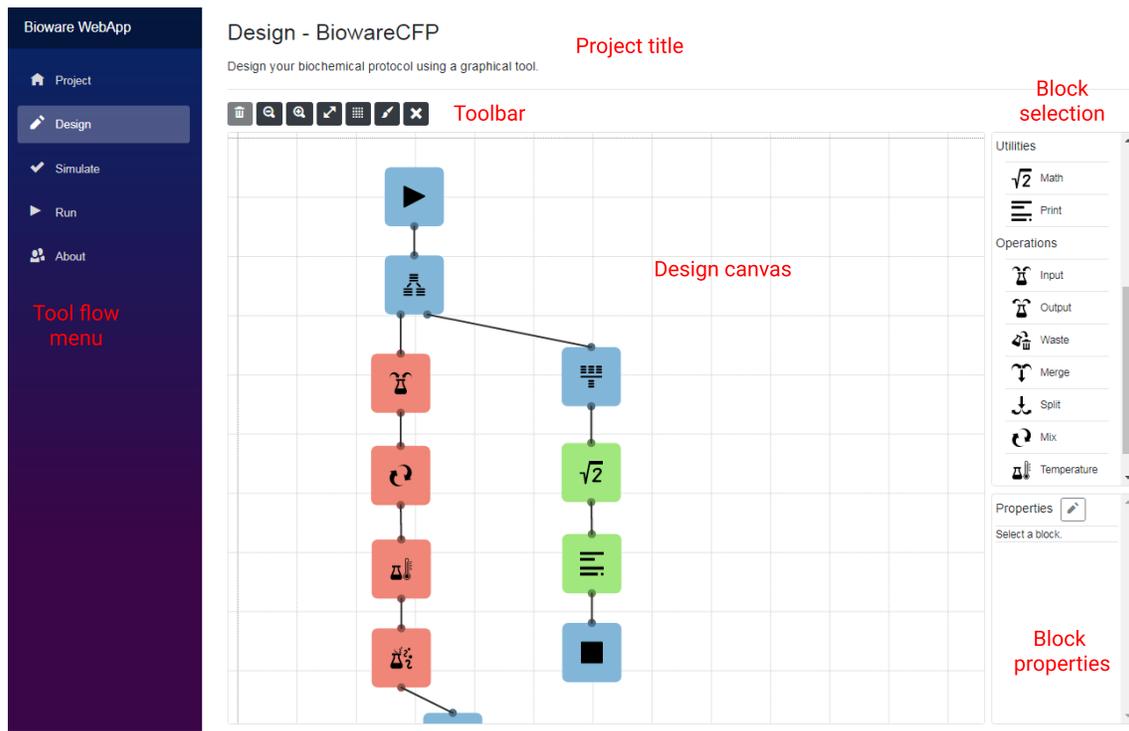


Figure 3: The web-based interface for designing a protocol using an intuitive high-level block-based representation. Blocks correspond to classic programming operations (arithmetics, flow control, printing, etc.) and fluidic operations (dispensing, mixing, incubation, temperature control, sensing, etc.).

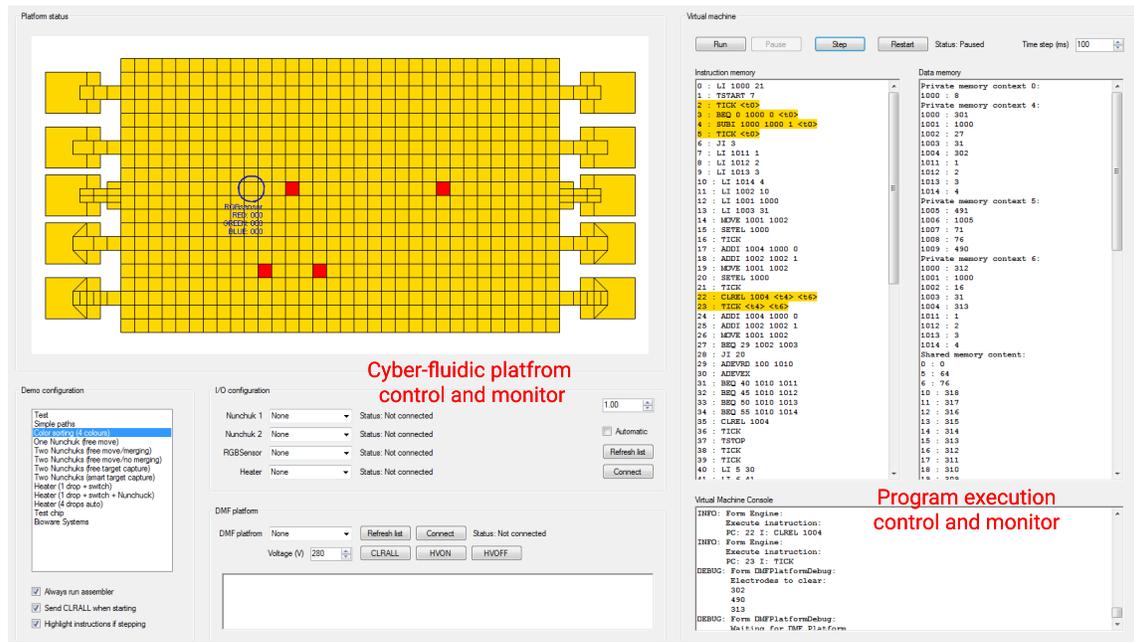


Figure 4: The graphical user interface of the virtual machine executing a testing protocol where droplets are sorted based on their color. The interface allows users to monitor and control the protocol execution and the cyber-fluidic platform.

REFERENCES

- [1] ALISTAR, M., AND GAUDENZ, U. OpenDrop: An integrated do-it-yourself platform for personal use of biochips. *Journal of Bioengineering* 4, 4 (2017), 45.
- [2] AMIN, A. M., THOTTETHODI, M., VIJAYKUMAR, T. N., WERELEY, S., AND JACOBSON, S. C. AquaCore: A programmable architecture for microfluidics. In *Proc. of the 34th Annual International Symposium on Computer Architecture* (2007), ACM, pp. 254–265.
- [3] ANANTHANARAYANAN, V., AND THIES, W. BioCoder: A programming language for standardizing and automating biology protocols. *Journal of Biological Engineering* 4, 1 (2010), 13.
- [4] FOBEL, R., FOBEL, C., AND WHEELER, A. R. DropBot: An open-source digital microfluidic control system with precise control of electrostatic driving force and instantaneous drop velocity measurement. *Applied Physics Letters* 102, 19 (2013), 1129–1132.
- [5] OTT, J., LOVELESS, T., CURTIS, C., LESANI, M., AND BRISK, P. BioScript: Programming safe chemistry on laboratories-on-a-chip. *Proc. of the ACM on Programming Languages* 2 (2018), 1–31.
- [6] POP, P., ALISTAR, M., STUART, E., AND MADSEN, J. *Fault-tolerant Digital Microfluidic Biochips: Compilation and Synthesis*. Springer, 2015.
- [7] SADOWSKI, M. I., GRANT, C., AND FELL, T. S. Harnessing QbD, programming languages, and automation for reproducible biology. *Trends in Biotechnology* 34, 3 (2016), 214–227.
- [8] TANEV, G., PEZZAROSSA, L., SVENDSEN, W., AND MADSEN, J. A modular reconfigurable digital microfluidics platform. In *Proc. of the 19th Symposium on Design, Test, Integration & Packaging of MEMS and MOEMS* (2018), IEEE, pp. 1–6.
- [9] TANEV, G., SVENDSEN, W., AND MADSEN, J. A reconfigurable digital microfluidics platform. In *Proc. of the 11th International Workshop on Bio-Design Automation* (2019), ACM, pp. 1–6.
- [10] TANEV, G. P. *A Modular Design Approach For Programmable Cyber-Fluidic Systems*. PhD thesis, Technical University of Denmark, 2021.
- [11] WILLSEY, M., STEPHENSON, A. P., TAKAHASHI, C., VAID, P., NGUYEN, B. H., PISZCZEK, M., BETTS, C., NEWMAN, S., JOSHI, S., STRAUSS, K., AND CEZE, L. Puddle: A dynamic, error-correcting, full-stack microfluidics platform. In *Proc. of the 24th International Conference on Architectural Support for Programming Languages and Operating Systems* (2019), pp. 183–197.

Comparison of Extrinsic and Intrinsic Noise Model Predictions for Genetic Circuit Failures

Pedro Fontanarrosa
pfontanarrosa@gmail.com
University of Utah
Salt Lake City, UT, USA

Lukas Buecherl
lukas.buecherl@colorado.edu
University of Colorado, Boulder
Boulder, CO, USA

Chris J. Myers
chris.Myers@colorado.edu
University of Colorado, Boulder
Boulder, CO, USA

1 INTRODUCTION

In recent years, *Genetic Design Automation* (GDA) tools like iBioSim [8, 18], SBOLCanvas [16], and Cello [9] were established to support scientists in automatically designing and modeling genetic circuits. This allows researchers to simulate their designs *in-silico* before building them in *in-vitro*, saving time and money. Therefore, mathematical descriptions of genetic networks become a necessity and that is why genetic design is usually model-driven [11]. The deterministic framework of ODE analysis is appropriate to describe the mean behavior of a system. However, there is no randomness or stochasticity associated with the model, and the same results are obtained given the same initial conditions [1]. The stochastic nature of biochemical reactions, even at the single-gene level [4], and the fluctuations in reaction rates due to differences in the environment [2, 10], generate significant noise to a system [7, 10, 12, 15]. Therefore, stochastic modeling and analysis is necessary to accurately predict outcome production fluctuations of GRNs. These fluctuations, caused by noise within a system, can have a drastic effect on the probabilities of glitching behavior of a system. This makes stochastic analysis of GRNs critical for the study of circuit failure probabilities, since failures can have drastic effects if the circuit’s output results in an early drug release or apoptosis.

Work by Fontanarrosa et al. [5] and Buecherl et al. [3] investigate the glitching behavior of a combinational genetic circuit first published by Nielsen et al. [9]. The genetic circuit was labeled 0x8E and can be seen in Figure 1(a). Fontanarrosa et al. [5] identified the different input transitions of the circuit that resulted in a glitching behavior. Based on this information, Fontanarrosa et al. designed two additional implementations of the circuit with the same function, using different logic combinations. The two modified layouts can be seen in Figure 1(b) and (c). The layout in (b) has redundant logic as two NOT-gates, to add delay to the IPTG pathway. The layout in (c) was created by using hazard-preserving optimization methods to avoid introducing *logic hazards*. Buecherl et al. [3] built on those results using stochastic simulation and stochastic model checking to determine the probability of the glitches discovered by Fontanarrosa et al.

There are different sources of noise that would generate variability in a circuit’s output. The inherent stochasticity of biochemical processes, such as transcription and translation, generates *intrinsic* noise [15]. This is especially significant in systems with low copy numbers of mRNAs or proteins in living systems [15, 17]. Therefore, stochastic effects are thought to be particularly important for gene expression and have been invoked to explain cell-cell variations of output production in clonal populations [4, 15]. The “stochastic chemical kinetics” that arise due to random births and deaths of individual molecules give rise to jump Markov processes, which can be analyzed by means of master equations and simulated with stochastic simulation algorithms [6, 7]. However, Beal [2] argues that stochastic chemical kinetics cannot explain the observed variation, and thus the explanation of such variation falls back to extrinsic noise. *Extrinsic* noise is generally defined as fluctuations and variability in a system’s reaction rates due to disturbances originated from its environment [10, 13]. This can be modeled as fluctuations in model parameters (such as transcription and degradation rates) [14].

This work focuses on determining if there are differences in predicted circuit failure percentages for three different circuit layouts with identical expected functions, using stochastic analysis to simulate different noise sources. The results shed light on the difference between the intrinsic and extrinsic noise model predictions and if the differences in circuit layouts have any effect on glitch propensities. This, in turn, emphasizes the need to evaluate further the relative influence of intrinsic and extrinsic noise on a genetic circuit’s output to help designers predict circuit failures more accurately and, therefore, determine better design choices. Moreover, the percent failure predictions between different circuit layouts can help designers weigh different options of circuit topologies to determine which one is best suited for the intended purposes of the design.

2 CIRCUIT ANALYSIS

Intrinsic Noise Model

In this work, to reduce the intrinsic model’s complexity, protein production and degradation was set to steps in molecule change to ten instead of one. That means that every time a

production or degradation reaction is fired, ten molecules are produced or destroyed with a ten times reduction in reaction propensity. Furthermore, to reduce the number of species, only the internal molecules and complexes are modeled instead of the input molecules *IPTG*, *aTc*, *Ara*. For instance, *IPTG* binds to *LacI* which regulates the circuit. Instead of modeling both species, only *LacI* is modeled.

Extrinsic Noise Model

The extrinsic noise model used in this work applies a simple case of *static* external perturbations, modeled as a random draw from a folded *normal* distribution for each parameter value used in the model at the beginning of each simulation run. The mean of each distribution is the default parameter value in iBioSim (obtained from literature), with a standard deviation equal to forty percent of the mean's absolute value (which emulates the "extrinsic noise"). This value of noise was obtained when calibrating different noise values, but is an arbitrary value that should be replaced with a better estimate obtained from experimentation (see Section 3). Beal [2] argues that these parameters follow a geometric distribution instead of a normal one, which is also part of our planned future work.

Results

For the simulation of the circuit models, the input molecules for the high state are set to 60 molecules. Any further increase of the molecule count does not affect the output of the circuit. The models use iBioSim's default parameters. The thresholds for static-0 function hazards and static-1 function hazards are defined by analyzing the response of a standardized NOT gate, which is ten and thirty molecules, respectively. More information can be found in [3]. The results of the intrinsic and extrinsic analysis can be seen in Table 1.

Table 1 shows the predicted circuit failure percentages for both the intrinsic and extrinsic analysis for the three different circuit layouts: Original Design (OD, Figure 1(a)), Redundant Logic (RL, Figure 1(b)), and Logic-Hazard-Free (LHF, Figure 1(c)). The results show that depending on which noise model is taken into consideration, the circuit layout with the lowest probability of glitching differs. So, for example, if only intrinsic noise model predictions are considered, then the LHF circuit layout is the best choice. Otherwise, if only extrinsic noise model predictions are considered, then the OD layout appears to be the best. Furthermore, the results show that the glitch propensities for each transition vary greatly depending on which noise model is considered for each circuit layout. For example, for input transition (1,1,1) to (0,1,0), 91% of the intrinsic noise simulation runs for the OD layout circuits glitch, whereas only 36% of the extrinsic noise simulation runs glitch for the same circuits.

3 DISCUSSION

This paper presents a comparison of genetic circuit failure percentage predictions from different models used to simulate noise in a GRN. It illustrates that the design choice is affected if noise is considered arising primarily from intrinsic or extrinsic sources, since the probabilities of circuit failures change for each simulation. In this case, intrinsic noise modeling generally predicts a higher percentage of glitching behavior for all circuits than extrinsic noise modeling. This motivates further study to determine which source of noise has a higher incidence in GRNs, if not both, to be able to accurately predict glitch propensities for genetic circuits. One current limitation of the work is that these genetic circuits operate in a range of 10^4 to 10^5 molecules for a high signal. In a future iteration of the work these aspects should be considered to reflect more biologically realistic values.

This paper also presents a comparison of genetic circuit failure percentage predictions for three different circuit layouts with identical expected function. If intrinsic noise has a higher influence on a circuit's output, then these results show that the LHF circuit layout is the least likely to glitch. Whereas if extrinsic noise is the predominant source of fluctuations in GRNs, then the OD design layout is the least likely to glitch.

In the future, we plan to extend the analysis done in this paper to include:

- characterized gate parameters for increased accuracy of results,
- percent circuit failures for all transitions, not only those that have function hazards,
- stability of each circuit for each state, calculated as the ability to hold-state given a fixed input state,
- extend the extrinsic noise model to use the model proposed by Beal [2],
- generate a model that simulates both intrinsic and extrinsic sources of noise.

These results and further studies will advance the DBTL pipeline, promoting the learning and designing stage to filter out the circuit layouts with higher probability of glitching for the input transitions considered critical to the designer. Furthermore, a model generator implemented in iBiosim that automatically includes intrinsic or extrinsic sources of noise in the model would help genetic circuit designers apply and test different levels of noise to obtain both circuit failure predictions and circuit robustness.

Acknowledgements

The authors of this work are supported by DARPA FA8750-17-C-0229 and by National Science Foundation Grant No.

Table 1: Intrinsic and extrinsic noise model predictions of circuit failure percentages

Input Transition	Intrinsic			Extrinsic		
	OD	RL	LHF	OD	RL	LHF
(0, 1, 0) → (1, 1, 1)	0.31	0.76	0.29	0.09	0.27	0.09
(0, 1, 0) → (1, 0, 0)	0.72	0.80	0.25	0.36	0.34	0.08
(1, 1, 1) → (1, 0, 0)	0.92	0.93	0.91	0.48	0.38	0.31
(1, 1, 1) → (0, 1, 0)	0.91	0.54	0.90	0.36	0.19	0.37
(1, 0, 0) → (0, 1, 0)	0.76	0.42	0.74	0.33	0.20	0.29
(1, 0, 0) → (1, 1, 1)	0.30	0.33	0.30	0.13	0.15	0.05
(0, 1, 1) → (1, 0, 1)	0.99	0.81	0.99	0.85	0.64	0.88
(0, 0, 0) → (0, 1, 1)	0.83	0.83	0.82	0.37	0.55	0.43
(0, 0, 0) → (1, 0, 1)	0.99	0.81	0.99	0.82	0.65	0.90
(1, 0, 1) → (0, 1, 1)	0.99	1.00	0.99	0.83	0.88	0.86
(0, 1, 1) → (0, 0, 0)	0.86	0.87	0.77	0.52	0.64	0.54
(1, 0, 1) → (0, 0, 0)	0.86	0.96	0.74	0.59	0.73	0.60

Percent failure predictions for each input transition that contains a function hazard, for different circuit layouts, and different noise models. The order of the inputs is *Ara*, *IPTG*, *aTc*. So, for example, (0,1,0) means only *IPTG* is present. OD: Original Design layout (Figure 1(a)), RL: Redundant Logic layout (Figure 1(b)), LHF: Logic-hazard free layout (Figure 1(c)).

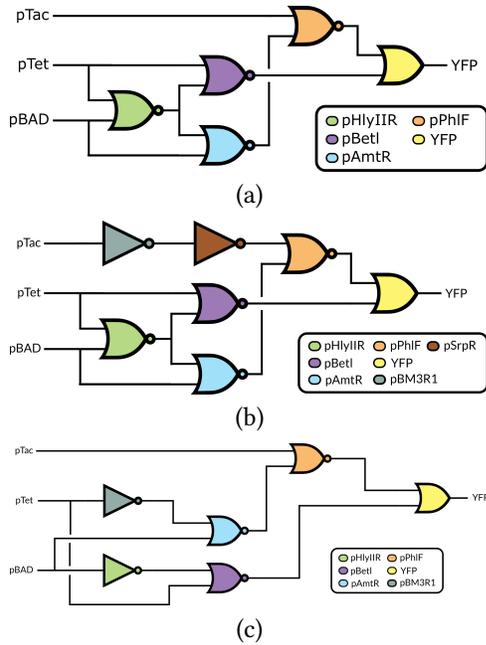


Figure 1: Three different logic layouts for the circuit 0x8E. The three inducer molecules are *IPTG*, *aTc* and *Ara* and the output is *YFP*. The OR gate is represented by \vee and the NOR gate by ∇ . (a) Original circuit layout as published in [9]. (b) Circuit implementation with added redundant logic as two NOT gates, which add an extra delay to the *IPTG* pathway. The NOT gate is represented by \neg . (c) Circuit implementation with logic-hazard-free optimizations. More details on the implementations (b) and (c) can be found in [5].

1856740. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the funding agencies.

REFERENCES

- [1] BALDWIN, G., BAYER, T., DICKINSON, R., ELLIS, T., FREEMONT, P. S., KITNEY, R. I., POLIZZI, K. M., AND STAN, G.-B., Eds. *Synthetic Biology: A Primer*, revised edition ed. Imperial College Press ; World Scientific Publishing Co. Pte. Ltd, [London] : Singapore ; Hackensack, NJ ; London, 2016.
- [2] BEAL, J. Biochemical complexity drives log-normal variation in genetic expression. *Engineering Biology* 1, 1 (June 2017), 55–60.
- [3] BUECHERL, L., FONTANARROSA, P., THOMAS, P. J., MANTE, J., ZHANG, Z., AND MYERS, C. J. Stochastic Hazard Analysis of Genetic Circuits in iBioSim and STAMINA. *ACS Synthetic Biology (Under Revision)* (2021).
- [4] ELOWITZ, M. B., LEVINE, A. J., SIGGIA, E. D., AND SWAIN, P. S. Stochastic gene expression in a single cell. *Science* 297, 5584 (Aug. 2002), 1183–1186.
- [5] FONTANARROSA, P., DOOSTHOSSEINI, H., ESPAH BORUJENI, A., DORFAN, Y., VOIGT, C. A., AND MYERS, C. J. Genetic Circuit Dynamics: Hazard and Glitch Analysis. *ACS Synthetic Biology* (Aug. 2020).
- [6] GILLESPIE, D. T. Exact stochastic simulation of coupled chemical reactions. *The Journal of Physical Chemistry* 81, 25 (Dec. 1977), 2340–2361.
- [7] LESTAS, I., PAULSSON, J., ROSS, N. E., AND VINNICOMBE, G. Noise in Gene Regulatory Networks. *IEEE Transactions on Automatic Control* 53, Special Issue (Jan. 2008), 189–200.
- [8] MYERS, C. J., BARKER, N., JONES, K., KUWAHARA, H., MADSEN, C., AND NGUYEN, N.-P. D. iBioSim: A tool for the analysis and design of genetic circuits. *Bioinformatics* 25, 21 (Nov. 2009), 2848–2849.
- [9] NIELSEN, A. A. K., DER, B. S., SHIN, J., VAIDYANATHAN, P., PARALANOV, V., STRYCHALSKI, E. A., ROSS, D., DENSMORE, D., AND VOIGT, C. A. Genetic circuit design automation. *Science* 352, 6281 (Apr. 2016), aac7341.
- [10] PAULSSON, J. Summing up the noise in gene networks. *Nature* 427,

- 6973 (Jan. 2004), 415–418.
- [11] PECCOUD, J. Synthetic biology: Fostering the cyber-biological revolution. *Synthetic Biology* 1, 1 (Jan. 2016).
- [12] PURNICK, P. E. M., AND WEISS, R. The second wave of synthetic biology: From modules to systems. *Nature Reviews Molecular Cell Biology* 10, 6 (June 2009), 410–422.
- [13] SCOTT, M., INGALLS, B., AND KERN, M. Estimations of intrinsic and extrinsic noise in models of nonlinear genetic networks. *Chaos: An Interdisciplinary Journal of Nonlinear Science* 16, 2 (June 2006), 026107.
- [14] SINGH, A., AND SOLTANI, M. Quantifying Intrinsic and Extrinsic Variability in Stochastic Gene Expression Models. *PLOS ONE* 8, 12 (Dec. 2013), e84301.
- [15] SWAIN, P. S., ELOWITZ, M. B., AND SIGGIA, E. D. Intrinsic and extrinsic contributions to stochasticity in gene expression. *Proceedings of the National Academy of Sciences* 99, 20 (Oct. 2002), 12795–12800.
- [16] TERRY, L., EARL, J., THAYER, S., BRIDGE, S., AND MYERS, C. J. Sbolcanvas: A visual editor for genetic designs. *ACS Synthetic Biology* (2021).
- [17] THATTAI, M., AND VAN OUDENAARDEN, A. Intrinsic noise in gene regulatory networks. *Proceedings of the National Academy of Sciences* 98, 15 (July 2001), 8614–8619.
- [18] WATANABE, L., NGUYEN, T., ZHANG, M., ZUNDEL, Z., ZHANG, Z., MADSEN, C., ROEHNER, N., AND MYERS, C. iBioSim 3: A tool for model-based genetic circuit design. *ACS Synthetic Biology* (June 2018).

Modeling of the engineered production of curcumin in *Escherichia coli*

Michael Cotner*
Ellie Brown*
michael.cotner@usu.edu
ellie.siddoway@usu.edu
Biological Engineering
Utah State University
Logan, UT, USA

Jixun Zhan
jixun.zhan@usu.edu
Biological Engineering
Utah State University
Logan, UT, USA

Zhen Zhang
zhen.zhang@usu.edu
Electrical and Computer Engineering
Utah State University
Logan, UT, USA

1 INTRODUCTION

Curcumin is the major bioactive natural product in turmeric (*Curcuma longa*), which is commonly used as a food additive (flavor and colorant) and traditional medicine for thousands of years. This polyphenol has a variety of biological activities, such as antioxidant[15], anti-cancer[18], anti-allergic[2], anti-inflammatory[13], and anti-Alzheimer's[7] effects. The production of this promising natural product relies on the extraction of producing plants, which requires large quantities of farmland and organic solvents[17]. Therefore, this traditional method is not environmentally friendly or cost-effective. Additionally, the quality of the product is often affected by the regions where the plants grow. Microbial production represents a great alternative because of the significantly reduced production time and reproducible production process. Microorganisms such as *Escherichia coli* can be engineered into curcumin-producing cell factories by incorporating curcumin biosynthetic enzymes. Our group has previously introduced five enzymes: tyrosine ammonia lyase (TAL), 4-Coumarate:CoA ligase (4CL), coumarate 3-hydroxylase (C3H), caffeic acid 3-O-methyltransferase (COMT) and curcuminoid synthase (CUS), into *E. coli* to construct an artificial pathway to produce curcumin in the bacterium.

To enable industrial production of curcumin, it is necessary to continuously improve the yield in engineered strains. While testing various parameters in lab experiments are often time-consuming and labor-intensive, this work aims to establish a computer model to simulate the production of curcumin in *E. coli*, based on which we could identify the limiting factors in the pathway for experimental optimization. We used the above-mentioned biosynthetic enzymes except that the CUS was replaced with two type III polyketide synthases, diketide-CoA synthase (DCS) and curcumin synthase (CURS) due to the lack of kinetic data for CUS.

2 PROBABILISTIC MODEL CONSTRUCTION

Curcumin production in *E. coli* was previously modeled to predict and optimize the yield[14]. Another model was created using ODEs to accurately predict the production[1]. In this work, we investigate the probabilistic behavior of the pathway model. The PRISM Probabilistic Model Checking language is used to model the pathway as a continuous-time Markov chain (CTMC)[12]. PRISM has previously been used to model the Fibroblast Growth Factor (FGF) pathway[8]. The pathway is shown in Figure 1, as modeled. Intracellular L-tyrosine serves as the starting substrate of this pathway, which is converted into feruloyl-CoA by TAL, C3H, 4CL, and COMT. Feruloyl-CoA is extended by DCS to produce feruloylacetyl-CoA, which is then condensed with another feruloyl-CoA molecule to produce curcumin through CURS. In the adaptation of the pathway to the probabilistic model, we made the following assumptions. First, the model abstracts away the pathway's interaction with the cell's native metabolism. The DCS enzyme utilizes malonyl-CoA as a substrate, a common intracellular species important to the cell's central metabolism and fatty acid biosynthesis. We approximate the concentration of malonyl-CoA as a constant assuming that it is always in sufficient supply for DCS to produce feruloylacetyl-CoA. Secondly, the model also excludes the byproducts of DCS and CURS when these enzymes act on caffeoyl-CoA and *p*-coumaroyl-CoA. In a laboratory system, curcuminoid byproducts of these enzymes, such as dicaffeoylmethane, demethoxycurcumin, and bisdemethoxycurcumin, would be produced at low concentrations along with curcumin. These byproducts were excluded due to a lack of their kinetic parameters in existing literature, but we hope to include them in the future. Our previous lab experiments indicate that these byproducts are only produced at trace concentrations [19]. Therefore, excluding them may not have a severe effect on the results of the model. The computational model for the pathway consists of nine species and ten reactions. Each species is represented by an integer variable, whose value corresponds to the micromolar intracellular concentration of that species. Each reaction is

*This research was funded by the Utah State University Engineering Undergraduate Research Program and gift donations from Adobe Incorporated.

modeled as a probabilistic transition that updates concentrations of participating species. The rate of each reaction is calculated using the Michaelis-Menten equation, a function of an enzyme’s molecular weight, k_{cat} and K_m for a specific enzyme and substrate, and concentration of the substrate. Table 1 shows the kinetic parameters of each enzyme and respective substrate found in literature. Each parameter was taken from enzymes from the organisms used in our group’s past resveratrol model[3]. Given the initial state of the model, which is determined by the initial concentration of each species, the reaction rates are calculated. Based on the ratio of the rates for all enabled reactions at the initial state, the model then probabilistically chooses a reaction to occur, resulting in decrementing its substrate concentration and incrementing the product concentration by one micromolar. In the next state, each reaction rate is reevaluated using the updated concentrations, and the process is repeated until a user-specified upper time bound is reached. To calculate the most probable, average behavior, the results of 1,000 simulations are averaged.

Table 1: Enzyme kinetics data

Enzyme	Substrate	$K_m(\mu\text{M})$	k_{cat} (1/s)	Source
TAL	L-tyrosine	1492.2	155	[4]
4CL	<i>p</i> -coumaric acid	25.1	16.3279	[6]
4CL	caffeic acid	11	4.408533	[5]
4CL	ferulic acid	56.7	7.475801	[6]
C3H	<i>p</i> -coumaric acid	143.03	0.0347	[9]
C3H	<i>p</i> -coumaroyl-CoA	143.03	0.0347	[9]
COMT	caffeic acid	59.5	0.145152	[16]
COMT	caffeoyl-CoA	26	0.145152	[10]
DCS	feruloyl-CoA	46	0.02	[11]
CURS	feruloyl-CoA	18	0.018333	[11]

Listing 1: Model of the curcumin production.

```

module partial_pathway
[] fcoa >0 -> v1:(fcoa' = fcoa + 1) & (fcoa' =
  fcoa - 1);
[] fcoa >0 & facoa >0 -> v2:(cur' = cur + 1) & (
  fcoa' = fcoa - 1) & (facoa' = facoa - 1);
endmodule

```

A model snippet is shown in Listing 1. The full model can be found at <https://github.com/formal-verification-research/phenylpropanoids-model>. Using feruloyl-CoA as a substrate, two reactions produce feruloyl-CoA and curcumin, respectively. Rate formulas v_1 and v_2 (not shown here) are evaluated using the concentration of feruloyl-CoA and each enzyme’s respective kinetic parameters at each simulation

step. These two reactions compete for feruloyl-CoA as a substrate and the ratio of their rates determines the likelihood of one reaction occurring over the other. As this example shows, the probabilistic behavior emerges from the random choice made between multiple competing reactions at each simulation step. Therefore, it facilitates a more in-depth analysis than a typical deterministic simulation, such as the likelihood of certain reactions can happen and limiting enzymes in the synthetic pathway.

3 RESULTS AND DISCUSSION

Figure 2 shows the results of an average of 100 simulations. It can be observed that every intermediate in the pathway eventually reaches an almost constant concentration except the product, curcumin, and *p*-coumaroyl-CoA. Curcumin reaches a steady linear increase, which is expected as there are no parts of the cell that use curcumin. The concentration of *p*-coumaroyl-CoA rapidly builds in the cell as 4CL, which produces the molecule, is much faster than C3H, which utilizes the molecule as a substrate. Both caffeic acid and ferulic acid are observed at an average concentration of 0 μM due to this same reason. While these two molecules are produced by the simulation, because 4CL is drastically more efficient than C3H, it is much more likely that *p*-coumaroyl-CoA is produced over caffeic acid. When caffeic acid is produced, however, it is quickly converted into caffeoyl-CoA by 4CL, again because of the enzyme’s high efficiency.

Based on these observations, C3H is likely the limiting enzyme and primary candidate for optimization. To validate this conclusion, an overexpression simulation was run on every enzyme in the pathway. In these simulations, the concentration of a selected enzyme was assumed to be double that of other enzymes. The concentrations of all other enzymes for each simulation were kept at the original value of 25 mg/L, and the overexpressed value was set to 50 mg/L. If the simulation saw an increase in yield, then these enzymes were marked as a potential point of optimization to increase the yield. The simulation results are shown in Table 2.

Table 2: Curcumin yields from overexpressed simulation

Enzyme	Curcumin yield (μM)	Percent change
-	486.565 \pm 0.751	-
TAL	517.835 \pm 2.742	6.427%
4CL	474.484 \pm 2.520	-2.483%
C3H	645.888 \pm 3.848	32.744%
COMT	488.752 \pm 2.290	4.495%
DCS	434.854 \pm 3.348	-10.628%
CURS	490.037 \pm 2.479	7.136%

Based on these experiments, overexpressing most enzymes had a moderately positive or negative effect on the yield of curcumin. C3H, as expected, was the only enzyme that showed a significant increase in yield, and in our future optimization experiments we will focus on this enzyme. Interestingly, the overexpression of DCS caused a 10% decrease in yield, likely because this enzyme takes available feruloyl-CoA away from CURS and prevents curcumin from forming. This marks DCS as another potential point of optimization, as it may be beneficial to lower the concentration of this enzyme to improve the yield. Furthermore, we will test how adjusting the concentrations of multiple enzymes affects curcumin production. Using the information from our representative model, we will validate the model in the lab.

REFERENCES

- [1] BOADA, Y., VIGNONI, A., PICÓ, J., AND CARBONELL, P. Extended metabolic biosensor design for dynamic pathway regulation of cell factories. *iScience* 23, 7 (2020), 101305.
- [2] CHANIYILPARAMPU, R. N., NAIR, A. K., PARTHASARATHY, K., GOKARAJU, G. R., GOKARAJU, R. R., BHUPHATIRAJU, K., MANDAPATI, V. N. S. R. R., AND SOMASHEKARA, N. Curcuminoids and its metabolites for the application in allergic ocular/nasal conditions, 2014. US Patent 20120010297A1.
- [3] COTNER, M., ZHAN, J., AND ZHANG, Z. A computational metabolic model for engineered production of resveratrol in *Escherichia coli*. *ACS Synthetic Biology* 10, 8 (2021), 1992–2001. PMID: 34237218.
- [4] CUI, P., ZHONG, W., QIN, Y., TAO, F., WANG, W., AND ZHAN, J. Characterization of two new aromatic amino acid lyases from actinomycetes for highly efficient production of *p*-coumaric acid. *Bioprocess and Biosystems Engineering* 43 (Mar 2020), 1287–1298.
- [5] EHLTING, J., BUTTNER, D., WANG, Q., DOUGLAS, C. J., SOMSSICH, I. E., AND KOMBRINK, E. Three 4-coumarate:coenzyme A ligases in *Arabidopsis thaliana* represent two evolutionarily divergent classes in angiosperms. *The Plant Journal* 19 (2002), 9–20.
- [6] EHLTING, J., SHINE, J. J. K., AND DOUGLAS, C. J. Identification of 4-coumarate:coenzyme A ligase (4CL) substrate recognition domains. *The Plant Journal* 27 (2001), 455–465.
- [7] GOOZEE, K. G., SHAH, T. M., SOHRABI, H. R., RAINEY-SMITH, S. R., BROWN, B., VERDILE, G., AND MARTINS, R. N. Examining the potential clinical value of curcumin in the prevention and diagnosis of Alzheimer’s disease. *British Journal of Nutrition* 115, 3 (2016), 449–465.
- [8] HEATH, J., KWIATKOWSKA, M., NORMAN, G., PARKER, D., AND TYMCHYSHYN, O. Probabilistic model checking of complex biological pathways. 32–47.
- [9] HEO, K. T., KANG, S.-Y., JANG, D. J.-H., AND HONG, D. Y.-S. Sam5, a coumarate 3-hydroxylase from *Saccharothrix espanaensis*: new insight into the piceatannol production as a resveratrol 3'-hydroxylase. *ChemistrySelect* 2 (2017), 8785–8789.
- [10] INOUE, K., PARVATHI, K., AND DIXON, R. A. Substrate preferences of caffeic acid/5-hydroxyferulic acid 3/5-O-methyltransferases in developing stems of alfalfa (*Medicago sativa* L.). *Archives of Biochemistry and Biophysics* 375 (2000), 175–182.
- [11] KATSUYAMA, Y., KITA, T., FUNA, N., AND HORINOUCHE, S. Curcuminoid biosynthesis by two type III polyketide synthases in the herb *Curcuma longa*. *Journal of Biological Chemistry* 284 (2009), 11160–11170.
- [12] KWIATKOWSKA, M., NORMAN, G., AND PARKER, D. PRISM 4.0: Verification of probabilistic real-time systems. In *Proc. 23rd International Conference on Computer Aided Verification (CAV'11)* (2011), G. Gopalakrishnan and S. Qadeer, Eds., vol. 6806 of *LNCIS*, Springer, pp. 585–591.
- [13] LESTARI, M. L., AND INDRAYANTO, G. Chapter three - curcumin. 113–204.
- [14] MACHADO, D., RODRIGUES, L. R., AND ROCHA, I. A kinetic model for curcumin production in *Escherichia coli*. *Biosystems* 125 (2014), 16–21.
- [15] MENON, V. P., AND SUDHEER, A. R. *ANTIOXIDANT AND ANTI-INFLAMMATORY PROPERTIES OF CURCUMIN*. Springer US, Boston, MA, 2007, pp. 105–125.
- [16] PARVATHI, K., CHEN, F., GUO, D., BLOUNT, J. W., AND DIXON, R. A. Substrate preferences of O-methyltransferases in alfalfa suggest new pathways for 3-O-methylation of monolignols. *The Plant Journal* 25 (2008), 193–202.
- [17] PAWAR, H., GAVASANE, A., AND CHOUDHARY, P. A novel and simple approach for extraction and isolation of curcuminoids from turmeric rhizomes. *Advances in Recycling & Waste Management* 06 (01 2018).
- [18] VERA-RAMIREZ, L., PÉREZ-LOPEZ, P., VARELA-LOPEZ, A., RAMIREZ-TORTOSA, M., BATTINO, M., AND QUILES, J. L. Curcumin and liver disease. *BioFactors* 39, 1 (2013), 88–100.
- [19] WANG, S., ZHANG, S., XIAO, A., RASMUSSEN, M., SKIDMORE, C., AND ZHAN, J. Metabolic engineering of *Escherichia coli* for the biosynthesis of various phenylpropanoid derivatives. *Metabolic Engineering* 29 (2015), 153–159.

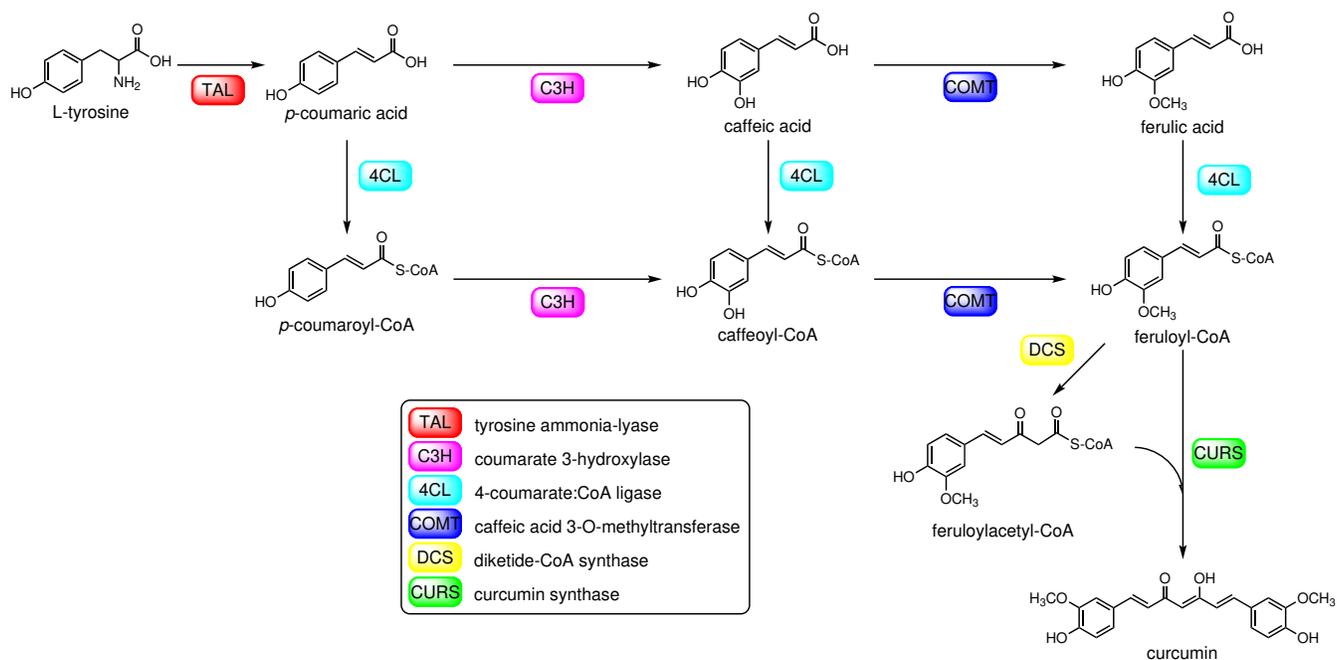


Figure 1: The modeled pathway

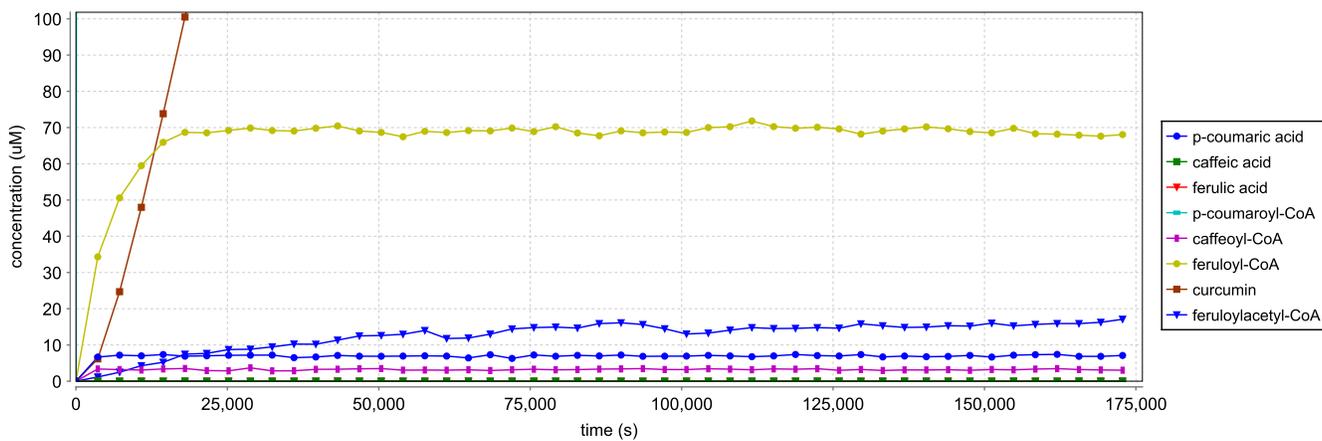


Figure 2: An average of 100 simulations

Biophysical Technology Mapping of Genetic Circuits

Nicolai Engelmann*, Tobias Schladt*, Erik Kubaczka*, Christian Hochberger, Heinz Koepl

TU Darmstadt, Germany

{nicolai.engelmann,heinz.koepl}@bcs.tu-darmstadt.de

{schlادت,hochberger}@rs.tu-darmstadt.de,erik.kubaczka@stud.tu-darmstadt.de

1 INTRODUCTION

Being still considered in its early stage of development compared to electronic design automation (EDA), genetic design automation (GDA) struggles to implement simple functions and accurately predict their behavior in the designated host environments. In terms of logic circuits, insufficient characterization of logic elements and the interaction with their host environment renders predictions of signals in the circuits imprecise. This is partially due to purely phenomenological approximations based on measurements of isolated compounds. This leads to strong requirements, like orthogonality, often limiting the number of available parts. On the level of DNA, transcription factors (TF's) encoded in genes, and promoters – which the TF's act on – are often chosen as the main building blocks implementing and connecting logic elements [7]. Employing available biophysical knowledge on this level, our design approach uses equilibrium thermodynamic interactions, thoroughly studied in recent works [1–3, 8], which takes competitive interactions, like crosstalk, and small copy-number effects into account. Selecting a candidate circuit in consideration of these effects in a technology mapping process, we ultimately seek to improve circuit performance and predictability in a real application. This, however, leads to longer evaluation times for each candidate demanding a clever search strategy. Therefore, we additionally propose a Branch and Bound based technology mapping scheme to obtain the optimal circuit while drastically reducing the number of candidate evaluations.

2 THERMODYNAMIC LOGIC CIRCUIT MODEL

Effects from small molecule copy-numbers and competitive binding like crosstalk and host genome interactions can alter a genetic circuit's response significantly up to orders of magnitude [3]. To be able to calculate the contributions of these effects, we need to consider sub-gate genetic parts but maintain sufficient scalability for a technology mapping process. The thermodynamic formulation provides a physical justification and economical parametrization but needs part specifications on the level of promoters, genes and TF's. However, for the calculations to remain transparent for higher layers, we continue to only allow the exchange of groups of the used sub-gate parts associated with the single logic

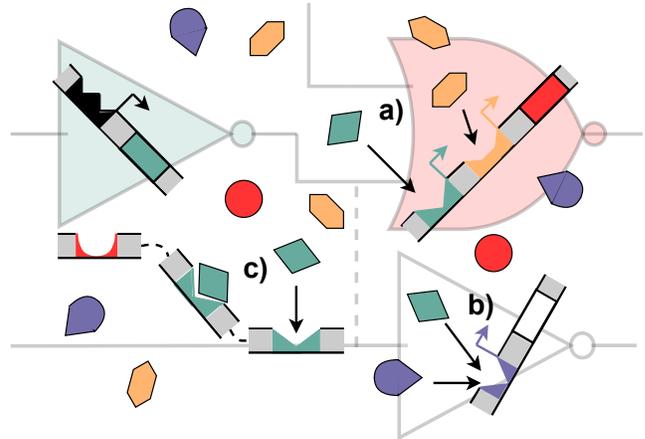


Figure 1: Illustration of the logic circuit from the thermodynamic perspective. Filled shapes expressed by the colour-coded genes bind to fitting promoters. The faint logic circuit in the background depicts the technology mapping process' point of view. Marked interactions are: a) desired binding implementing the circuit function, b) crosstalk, c) environmental competing binding sites.

gates. As our gate model we use the combination of a gene expressing a TF which represses a set of associated output promoters, inherited from the prominent GDA tool Cello [7]. The thermodynamic perspective of such a logic circuit is illustrated in Fig. 1.

Thermodynamic calculations on a candidate assignment are founded on the idea of time-scale separation of RNA-polymerase- (RNAP-) and TF-binding to available binding sites from protein-level signals which generate the circuit's response [2]. Combined with the assumption that expectation of RNAP-binding, and gene expression are proportional [8], we can calculate the circuit response from equilibrium statistics of RNAP- and TF-binding. Summations of statistical weights of microstates associated with specific events (i.e. partition functions) normalized to the total weight of all microstates give the expression levels. Thus, let $1, \dots, N$ be an enumeration of all gates and their associated parts. Let $\mathbf{x} \equiv (x_1, \dots, x_K) \in \mathbb{R}^K$ be a vector of all $K > N$ TF copy-numbers in the circuit with the first N being the gates' outputs followed by those posing as inputs and those from the host interfering with the circuit. Let further $Z(\mathbf{x})$ be the

* The authors contributed equally to this work.

total partition function of all microstates $s \in \mathcal{S}(x)$ of RNAP- and TF-binding permutations (we assume the number of RNAP to be fixed) and $Z_n(x)$ be the partition function of all microstates $s_n \in \mathcal{S}_n(x) \subset \mathcal{S}(x)$ where RNAP is bound to one or more input promoters upstream of gene n . Note, that $Z_n(x)$ is dependent on the circuit topology, since the set of input promoters of gene n can vary. We can find

$$x_n \propto \frac{Z_n(\mathbf{x})}{Z(\mathbf{x})} = \frac{\sum_{s_n \in \mathcal{S}_n(\mathbf{x}) \subset \mathcal{S}(\mathbf{x})} \exp(-\beta \varepsilon(s_n))}{\sum_{s \in \mathcal{S}(\mathbf{x})} \exp(-\beta \varepsilon(s))} = E(P_n), \quad (1)$$

where $\beta \equiv (k_B T)^{-1}$ is a product of the Boltzmann constant k_B and the equilibrium temperature T of the mix. The $\varepsilon(s)$ is the total binding energy for the permutation encoded in s . As an example, we provide the description of a single L -input NOR gate with inter-input crosstalk, such from $K - L$ environmental TF's, and leakage. There, we obtain from (1)

$$x_n \approx b_n \sum_{l=1}^L \frac{\exp\left(-\beta \left(\varepsilon_p^{(l)} - \varepsilon_{c,p}\right)\right)}{\sum_{i=1}^L \exp\left(-\beta \left(\varepsilon_p^{(i)} - \varepsilon_{c,p}\right)\right)} d^{(l)}(\mathbf{x}), \quad (2)$$

where the individual promoter-specific regulation factors $d^{(l)}$ are given by

$$d^{(l)}(\mathbf{x}) = \frac{1 + \sum_{k=1}^K \frac{x_k}{C} \exp\left(-\beta \left(\varepsilon_{p,x_k}^{(l)} - \varepsilon_p^{(l)} - \varepsilon_{c,x_k}\right)\right)}{1 + \sum_{k=1}^K \frac{x_k}{C} \exp\left(-\beta \left(\varepsilon_{x_k}^{(l)} - \varepsilon_{c,x_k}\right)\right)},$$

where we introduce the unrepressed copy-number $x_n \equiv b_n$ if $\mathbf{x} = \mathbf{0}$, the binding energies of RNAP to the L promoters $\varepsilon_p^{(l)}$, those to the C "background" sites on the host genome $\varepsilon_{c,p}^{(l)}$, the similar binding energies of the K TF's $\varepsilon_{f_k}^{(l)}$ and ε_{c,f_k} , the binding energy modelling imperfect competition between RNAP and TF $\varepsilon_{p,f_k}^{(l)}$ [10], and resolved the proportionality. Coming back to the full circuit with its N gates, we generally have the relations (c.f. (2))

$$\forall n : \quad x_n = b_n \frac{Z_n(\mathbf{x})Z(\mathbf{0})}{Z(\mathbf{x})Z_n(\mathbf{0})}, \quad (3)$$

with $\mathbf{0}$ the zero vector. Solving the implicit non-linear system (3) of N unknowns in N equations as a root-finding problem, we obtain the N unresolved quantities in \mathbf{x} using a Quasi-Newton algorithm. Equipped with this, we estimate crosstalk in the overlaid logic circuit, can incorporate competitive interaction with the host genome [3], and can potentially account for titration effects [8] in low copy-number regimes of involved TF's.

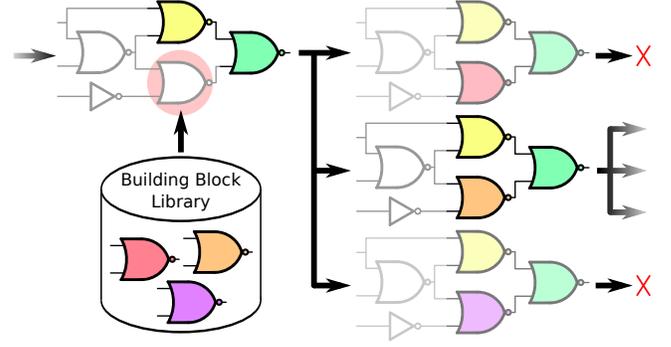


Figure 2: Excerpt of Branch and Bound's search tree for an exemplary circuit. The marked gate is systematically assigned library gates and resulting partial mappings are scored with the bounding function, i.e. simulated with optimistically estimated inputs. Inferior partial solutions are then pruned, while here, one solution is expanded further.

3 OPTIMAL AND FAST TECHNOLOGY MAPPING USING BRANCH AND BOUND

During genetic technology mapping, genetic building blocks from a library are assigned to the abstract logic circuit structure, while the resulting genetic circuit maximizes a given performance metric represented by a circuit score like output fold-change.

Due to the small circuit and library sizes and simple gate models, it has been feasible in GDA to find optimal technology mapping solutions using exhaustive search. However, increased model complexity, like it is introduced by the proposed thermodynamic model, leads to a computationally costly evaluation of candidate circuits. To handle this, heuristic approaches like Simulated Annealing have been proposed, which provide shorter run times, but do not deterministically obtain the optimal solution [7] [9]. As the output fold-change of genetic circuits heavily depends on the performance of the genetic building blocks near the end of the logic cascade, a Branch and Bound (B&B) based technology mapping scheme is proposed during which circuits are built iteratively starting at the primary output [4] [6]. This approach seeks to combine the computational performance of heuristics and the optimal solution quality of exhaustive search.

During B&B, an implicit enumeration of the search tree is performed, while based on the preliminarily best known solution only partial solutions with a prospect of leading to better solutions are considered further. As a result, the possibility of pruning parts of the search tree arises (see Fig. 2). The quality of partial solutions is given by the maximally reachable score of circuits containing the partial solutions, which is called the bounding function. It is based on a true partial simulation of the circuit, incorporating optimistic

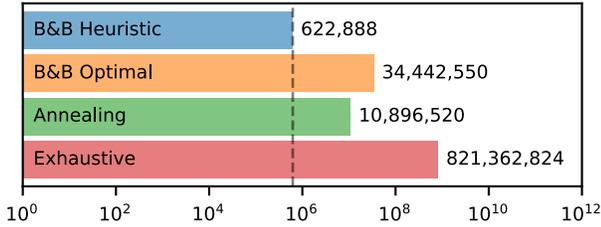


Figure 3: Number of simulations needed for mapping the set of 66 benchmark circuits with the examined algorithms.

bounds for repressor concentrations yet unknown and omitting un-reachable internal logic states identified by Boolean simulation.

The performance of this new technology mapping approach has been compared to exhaustive enumeration and the mean performance of the SA proposed in [9]. To this end, the genetic gate library and the gate model of Cello have been used to map a total of 66 circuit structures including the 33 structures presented in [7]. Furthermore, two configurations of B&B have been evaluated: an optimal one which guarantees finding the optimal solution and a heuristic one with a relaxed bounding function. Figure 3 shows the total number of simulations required by each mapping algorithm for the whole benchmark set. Compared to exhaustive search, B&B reduces the number of simulations 24-fold, while also obtaining the optimal technology mapping. When it comes to an heuristic approach, B&B outperforms SA 17.5-fold in run time, while increasing the chance of finding the optimal result 1.2-fold and decreasing the mean loss in score 6.1×10^5 -fold, leading to near-optimal results.

4 CONCLUSION

We propose a method to calculate a genetic logic circuit which uses a NOR-logic structure in the style of Cello [7] on a sub-gate level using a thermodynamic model instead of using gate transfer functions derived from kinetic models. The model uses a different library of elementary parts compared to a gate library but can be integrated in the more abstract gate-level technology mapping process by restricting the exchange of the parts. This allows taking small copy-number effects and such from competitive binding, like crosstalk, into account when scoring a circuit assignment. Statistical quantities used to characterize the sub-gate units can be inferred from experimental data [5]. Since these calculations are more expensive, we additionally propose a novel technology mapping approach using Branch and Bound, which shows to be capable of drastically reducing the average number of genetic circuits to evaluate before reaching an optimal solution, while using the classical gate model of Cello.

This is accomplished by rigorously pruning the search space with a domain-specific bounding function. Further work targets the incorporation of crosstalk in the bounding function, yielding tighter bounds with respect to the thermodynamic calculations, thus enabling a fast crosstalk-aware technology mapping.

REFERENCES

- [1] BINTU, L., BUCHLER, N., GARCIA, H., GERLAND, U., HWA, T., KONDEV, J., KUHLMAN, T., AND PHILLIPS, R. Transcriptional regulation by the numbers: Applications. *Current Opinion in Genetics and Development* 15, 2 (Apr. 2005), 125–135.
- [2] BINTU, L., BUCHLER, N., GARCIA, H., GERLAND, U., HWA, T., KONDEV, J., AND PHILLIPS, R. Transcriptional regulation by the numbers: Models. *Current opinion in genetics & development* 15 (05 2005), 116–24.
- [3] BREWSTER, R., WEINERT, F. M., GARCIA, H. G., SONG, D., RYDENFELT, M., AND PHILLIPS, R. The transcription factor titration effect dictates level of gene expression. *Cell* 156 (2014), 1312–1323.
- [4] CLAUSEN, J. Branch and bound algorithms-principles and examples. *Department of Computer Science, University of Copenhagen* (1999), 1–30.
- [5] GARCIA, H., AND PHILLIPS, R. Quantitative dissection of the simple repression input-output function. *Proceedings of the National Academy of Sciences of the United States of America* 108 (07 2011), 12173–8.
- [6] MORRISON, D. R., JACOBSON, S. H., SAUPPE, J. J., AND SEWELL, E. C. Branch-and-bound algorithms: A survey of recent advances in searching, branching, and pruning. *Discrete Optimization* 19 (2016), 79–102.
- [7] NIELSEN, A. A. K., DER, B. S., SHIN, J., VAIDYANATHAN, P., PARALANOV, V., STRYCHALSKI, E. A., ROSS, D., DENSMORE, D., AND VOIGT, C. A. Genetic circuit design automation. *Science* 352, 6281 (2016).
- [8] RYDENFELT, M., COX III, R., GARCIA, H., AND PHILLIPS, R. Statistical mechanical model of coupled transcription from multiple promoters due to transcription factor titration. *Physical review. E, Statistical, nonlinear, and soft matter physics* 89 (01 2014), 012702.
- [9] SCHLADT, T., ENGELMANN, N., KUBACZKA, E., HOCHBERGER, C., AND KOEPL, H. Automated design of robust genetic circuits: Structural variants and parameter uncertainty. *bioRxiv* (2021). BIORXIV/2021/456094.
- [10] SWANK, Z., LAOHAKUNAKORN, N., AND MAERKL, S. Cell-free gene-regulatory network engineering with synthetic transcription factors. *Proceedings of the National Academy of Sciences* 116 (03 2019), 201816591.

Automated translation of logical models to SystemVerilog enables simulation speedup

Eric Li¹, Emilee Holtzapple¹, Niteesh Sundaram¹, Natasa Miskov-Zivanov¹

¹University of Pittsburgh
 {erl75,erh87,nis101,nmzivanov}@pitt.edu

1 INTRODUCTION

Our understanding of disease and normal cell function is greatly enhanced by computational models of signaling cascades. Computational modeling methods, in contrast to *in vitro* or *in vivo* methods, require less time to produce mechanistic explanations of cell function. These models contain individual elements, which are usually a mix of genes, proteins, small molecules, and biological processes. For all but the simplest signaling cascades, any computational model will be large and complex. Each model element is capable of having multiple interactions, and feedback and feedforward loops are frequent in signaling cascades [3].

The complexity of cell signaling necessitates simulation methods that can manage large models, while still providing accurate results. Boolean models, described in [2], represent the activity of model elements as a discrete value, determined by logical functions (Figure 1 A,B). This methodology confers several benefits over other modeling approaches, as it is not dependent on reaction rates. This abstraction reduces the time needed for model assembly, as well as allowing the inclusion of more ambiguous model elements, such as biological processes like “inflammation”.

Software simulators are often limited by the typical runtimes for the programming language in which they are written. This is an area where we can benefit from the use of hardware simulators and, more specifically, Field-Programmable Gate Arrays (FPGAs). FPGAs provide better spatial parallelism, instruction efficiency, and re-programmability, all of which lead to a significant speedup in runtime (compared to software implementations). Recently, several hardware-based implementations of a simulator were proposed [4, 7, 8], written using hardware description languages (HDLs), and allowing several orders of magnitude of speedup compared to the software-based simulator. HDLs allow for the specification of gate-level logic and the creation of configurable and scalable models. Currently, however, there are no flexible, automated methods to convert Boolean models to an HDL. Here, we present a tool for converting Boolean models to SystemVerilog [1], the HDL used in this work (Figure 1C), for simulation of the model in an FPGA. This translation tool requires less time and manual intervention to produce models capable of being simulated using hardware methods.

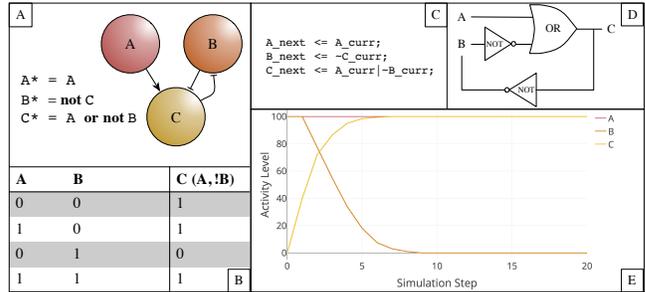


Figure 1: (A) A graphical representation of a simple toy Boolean model with three elements, as well as the Boolean expressions for the next state (*) for each element. (B) The truth table for the value of Element C. (C) SystemVerilog logic statements. (D) Logic gate representation of the toy model. (E) Simulation results using the DiSH simulator [9], a software-based simulator.

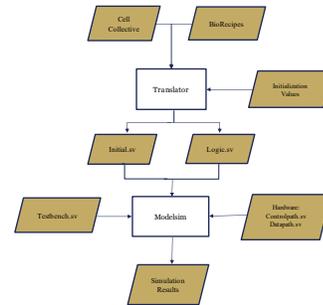


Figure 2: SystemVerilog translator workflow.

2 METHODS

In this work, a Python script was used to translate a logical model in the format of either BioRECIPIES or Boolean expressions into a hardware description language (HDL). The translated model is then used as input for a hardware simulator. The BioRECIPIES format is described in [10] and the Boolean expressions were found from models in Cell Collective. For each new model, the network logic module and initializations are changed and created by the translator. These are later combined with the controlpath and datapath modules from the hardware-based simulator, which remain the same from model to model (Figure 2).

To simulate each model, we use the hardware implementation of the Discrete Stochastic Heterogeneous Simulator (DiSH), described in detail in [4]. The DiSH simulator includes the datapath and controlpath modules, which provide ModelSim with the logic of how the update rules are used, namely which elements get updated and how frequently this occurs. In the hardware simulator, registers are the main storage components used to keep track of the network state. As such, the main components needed for the network logic module are the declaration of input and output registers for each element and the assign statement that links these registers. The translator handles the creation of these components on an element-by-element basis as the conditions of each element must be considered. In other words, the translator determines the regulators for each element from the input and, by reading the notation associated with each format, assigns the corresponding logic operator. When the entire model is translated to SystemVerilog and the model module is assembled, the model is simulated using ModelSim, a tool that is used to simulate standard HDL designs in a hardware environment. ModelSim compiles and links the individual SystemVerilog modules into a project where it proceeds to simulate the design. As the purpose of this work is to demonstrate the versatility of the translator in converting many different models, the simulator scheme and parameters were kept uniform and simple. The model is run under the simultaneous update scheme, which updates the activity of all model elements, based on their update rules, at each step in the simulation. Each model is run for 50 steps. While here we use the simultaneous update scheme, our SystemVerilog translator works with any DiSH update scheme.

3 RESULTS

To demonstrate the efficacy of the translator, we used as input three Boolean models Table 1 written in different representation formats. The runtime, even for large models with hundreds of nodes and edges, is in the order of milliseconds. Since the number of SystemVerilog logic statements relies only on the number of model elements, even the largest Boolean models can be translated quickly. For conversion of Boolean models to SystemVerilog, this is a profound speedup over manual methods. In addition to the time saving benefits of using the translator to convert large Boolean models, the translated models can be quickly simulated to explore model behavior. In the following section, we show the translation and simulation results for these Boolean models.

Large Granuloma Leukocyte Model

We translated a model of large granular lymphocyte leukemia (LGL) signal transduction, described in [11], from BioRECIPES format. Figure 3A is a graphical representation of the LGL model, accompanied by a close-up of a subnetwork. This part

Table 1: Characteristics of the test case Boolean models. For each model, we list the input format, number of elements in the model (nodes), and number of interactions within the model (edges). We also state the runtime of the translator in milliseconds (ms).

Model	Format	Elements (#)	Interactions (#)	Translation runtime (ms)
LGL	BioRecipes	47	168	9.94
ErbB	Cell Collective	247	1100	13.96
Fibroblasts	Cell Collective	139	557	6.98

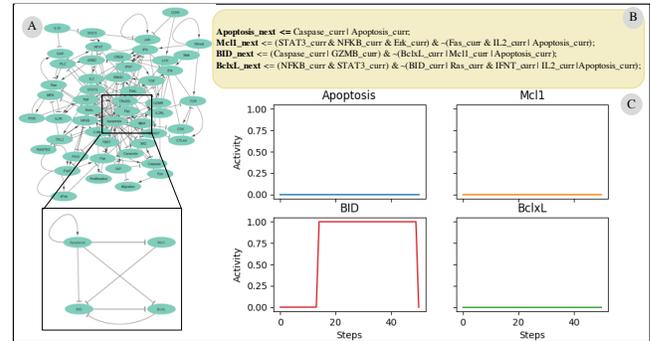


Figure 3: Translation of the LGL logical model described in [11]. (A) A graphical representation of the whole model, as well as a close-up of one pathway. The close-up shows several possible paths between the model elements Apoptosis and BclxL, both of which were focal points in the original publication. (B) A small subset of the SystemVerilog logic statements produced by the translator. We show the SystemVerilog statements for each model element in the close-up in (A). Note that some of the more verbose logic statements have been abridged for clarity. (C) Simulation results for the six elements in (B).

of the model represents the regulation of apoptosis by several elements - BclxL, Mcl1, and BID. The model is represented by 47 SystemVerilog logic statements after translation. The abridged SystemVerilog logic statements for the subnetwork are shown in Figure 3B. After translation of the LGL model into SystemVerilog, we simulated using the simultaneous update scheme, as described in the Methods section. The traces for the four subnetwork elements are shown in Figure 3C.

ErbB Signal Transduction Model

We translated a model of ErbB signaling, described in [5], to SystemVerilog. Figure 4A is a graphical representation of the model, with a close-up visualization of a subnetwork. This subnetwork represents EGF-dependent activation of Akt, which was studied in-depth in the original paper. This subnetwork contains six elements - EGF, EGFR dimer, modified EGFR dimer, free EGFR, PP2A, and Akt. The model, originally in Boolean expression format, is represented by 247 SystemVerilog logic statements after translation. The

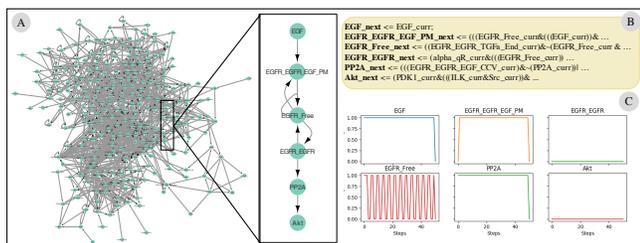


Figure 4: Translation of the ErbB logical model described in [5]. (A) A graphical representation of the whole model, as well as a close-up of one pathway. The close-up shows one possible path between the model elements EGF and Akt, both of which were focal points in the original publication. **(B)** A small subset of the SystemVerilog logic statements produced by the translator. We show the SystemVerilog statements for each model element in the close-up in (A). Note that some of the more verbose logic statements have been abridged for clarity. **(C)** Simulation results for the six elements in (B).

abridged SystemVerilog logic statements for the subnetwork are shown in Figure 4B. After translation of the LGL model into SystemVerilog, we simulated as described in Materials and Methods. The traces for the six subnetwork elements are shown in Figure 4C.

Fibroblast Model

We translated the Boolean model of fibroblast signaling transduction, described in [6]. Figure 4A is a graphical representation of the fibroblast model, accompanied by a close-up of one subnetwork. This part of the model represents the regulation of Erk through two pathways. The model, originally in Boolean expression format, is represented by 139 SystemVerilog logic statements after translation. The abridged SystemVerilog logic statements for the subnetwork are shown in Figure 4B. After translation of the LGL model into SystemVerilog. The traces for the six subnetwork elements are shown in Figure 4C.

4 CONCLUSION

Hardware-based simulators allow for orders of magnitude faster simulation than software simulators, but require input to be in the form of a hardware description language. Our translator is capable of fast conversion of Boolean models to SystemVerilog. Our results show that our tool translates even the largest models quickly, and the runtime in milliseconds does not neutralize the speedup of the DiSH hardware simulator over software simulation methods. This translator is also flexible in terms of the accepted input formats, and is extendable to non-biological models. We show in our three test cases that the translator automatically creates all

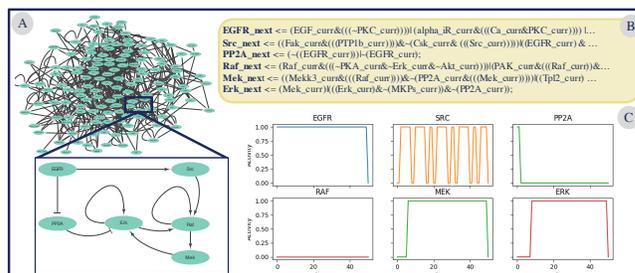


Figure 5: Translation of the fibroblast signaling transduction logical model described in [6]. (A) A graphical representation of the whole model, as well as a close-up of one pathway. The close-up shows several possible paths between the model elements EGFR and ERK, both of which were focal points in the original publication. **(B)** A small subset of the SystemVerilog logic statements produced by the translator. We show the SystemVerilog statements for each model element in the close-up in (A). Note that some of the more verbose logic statements have been abridged for clarity. **(C)** Simulation results for the six elements in (B).

SystemVerilog logic statements from the input model, and that these logic statements are still human readable. DiSH hardware simulation of these models demonstrates the utility of the translator tool, and its potential future uses for mechanistic study of cell signaling. Future directions include translation of other input formats, such as truth tables or adjacency matrices. The translator could also be updated to include discrete models with three or more activity levels, as opposed to Boolean models with only two activity levels. The DiSH software simulator is also capable of including temporal information, such as delays, within update rules. Potentially, we could include this functionality to the SystemVerilog translator.

Funding

Defense Advanced Research Projects Agency (W911NF-17-1-0135).

REFERENCES

- [1] Ieee standard for systemverilog—unified hardware design, specification, and verification language. *IEEE Std 1800-2017 (Revision of IEEE Std 1800-2012)* (2018), 1–1315.
- [2] ALBERT, I., THAKAR, J., LI, S., ZHANG, R., AND ALBERT, R. Boolean network simulations for life scientists. *Source Code Biol Med* 3 (2008), 16.
- [3] BRANDMAN, O., AND MEYER, T. Feedback loops shape cellular signals in space and time. *Science (New York, N.Y.)* 322, 5900 (2008), 390–395.
- [4] GILBOY, K., SAYED, K., SUNDARAM, N., BOCAN, K., AND MISKOV-ZIVANOV, N. A faster dish: Hardware implementation of a discrete cell signaling network simulator. In *2018 IEEE International Symposium on Circuits and Systems (ISCAS)*, pp. 1–5.

- [5] HELIKAR, T., KOCHI, N., KOWAL, B., DIMRI, M., NARAMURA, M., RAJA, S. M., BAND, V., BAND, H., AND ROGERS, J. A. A comprehensive, multi-scale dynamical model of erbb receptor signal transduction in human mammary epithelial cells. *PLOS ONE* 8, 4 (2013), e61757.
- [6] HELIKAR, T., KONVALINA, J., HEIDEL, J., AND ROGERS, J. A. Emergent decision-making in biological signal transduction networks. *Proceedings of the National Academy of Sciences* 105, 6 (2008), 1913–1918.
- [7] MISOV-ZIVANOV, N., BRESTICKER, A., KRISHNASWAMY, D., VENKATAKRISHNAN, S., KASHINKUNTI, P., MARCULESCU, D., AND FAEDER, J. R. Regulatory network analysis acceleration with reconfigurable hardware. In *2011 Annual International Conference of the IEEE Engineering in Medicine and Biology Society* (2011), pp. 149–152.
- [8] MISOV-ZIVANOV, N., BRESTICKER, A., KRISHNASWAMY, D., VENKATAKRISHNAN, S., MARCULESCU, D., AND FAEDER, J. R. Emulation of biological networks in reconfigurable hardware. In *Proceedings of the 2nd ACM Conference on Bioinformatics, Computational Biology and Biomedicine* (New York, NY, USA, 2011), BCB '11, Association for Computing Machinery, pp. 536–540.
- [9] SAYED, K., KUO, Y.-H., KULKARNI, A., AND MISOV-ZIVANOV, N. Dish simulator: Capturing dynamics of cellular signaling with heterogeneous knowledge. *2017 Winter Simulation Conference (WSC)*, pp. 896–907.
- [10] SAYED, K., TELMER, C. A., BUTCHY, A. A., AND MISOV-ZIVANOV, N. Recipes for translating big data machine reading to executable cellular signaling models. *Machine Learning, Optimization, and Big Data*, Springer International Publishing, pp. 1–15.
- [11] ZHANG, R., SHAH, M. V., YANG, J., NYLAND, S. B., LIU, X., YUN, J. K., ALBERT, R., AND LOUGHRAN, T. P. Network model of survival signaling in large granular lymphocyte leukemia. *Proceedings of the National Academy of Sciences* 105, 42 (2008), 16308–16313.

Excel-SBOL Converter: Creating SBOL from Excel Templates and Vice Versa

Jeanet Mante

jet@mante.net

University of Colorado Boulder

Isabel Pöttsch

imp31@cam.ac.uk

University of Cambridge

Julian Abam

julian.abam@colorado.edu

University of Colorado Boulder

Jacob Beal

jakebeal@ieee.org

Raytheon BBN Technologies

Chris J. Myers

chris.myers@colorado.edu

University of Colorado Boulder

1 INTRODUCTION

Synthetic biology is bringing together engineers and biologists [10]. Associated with this interdisciplinary movement is the need for reusable tools that supplement the current understanding of genetic sequences. To satisfy this need, Synthetic Biology communities across the world have developed tools and ontologies to help describe their unique semantic annotations [1, 3–9, 13, 14, 17–19, 22]. Shared representations for data and metadata, grounded in well-defined ontology terms, can help reduce confusion when sharing materials between practitioners.[20]. The Synthetic Biology Open Language (SBOL) [5] is one of the approaches that has been developed to address this challenge. SBOL provides a standardized format for the electronic exchange of information on the structural and functional aspect of biological designs, supporting use of engineering principles of abstraction, modularity, and standardization in synthetic biology. Many tools have been created that work with SBOL, including the SynBioHub repository software for storing and sharing designs [12].

Using formal representations such as SBOL, however, typically requires either a thorough understanding of these standards or a suite of tools developed in concurrence with the ontologies [11]. Unfortunately, this poses a significant barrier to use for scientists not trained to work with such abstractions. One approach to lowering this barrier was demonstrated in the Systems Biology for Micro-Organisms (SysMO) consortium [2]. In SysMO, the MicroArray Gene Expression Markup Language (Mage-ML) was set up as an XML schema [15], and users were expected to submit data to the SysMO Assets Catalogue (called SEEK) in XML format in order to publish work. To allow the use of the Mage-ML language without having to understand XML, the RightField tool was created [21], an ontology annotation and information management application that can add constrained ontology term selection to Excel spreadsheets. This tool enables administrators to create templates with controlled vocabularies, such that the scientists utilizing the tool would never actually see

the raw RightField, only the more familiar Excel spreadsheet interface.

Similarly, users of SBOL and SynBioHub have faced a steep learning curve for understanding the underlying ontology: as assessed in [16], “For successful use and interpretation of metadata presented in SynBioHub, the semantic annotation process should be biologist-friendly and hide the underlying RDF predicates.” Accordingly, the Excel-SBOL Converter presented here has been designed to provide a simple way for users to generate SBOL data without needing a detailed understanding of the underlying ontology and associated technologies. The converter provides a simple way for users to manage data by allowing users to download SBOL into Excel templates and submit Excel templates for conversion into SBOL.

2 RESULTS

The Excel-SBOL Converter enables the round trip conversion of SBOL2 into Excel Templates and Excel Templates to SBOL2, as shown in Figure 1. This metadata conversion allows researchers the flexibility of working in Excel without losing the benefits of SBOL2 and maintains inter-operability with the SBOL2 tool suite. The converter also facilitates understanding of the metadata relationships, management, and structure within the SBOL2 XML file. The converter has been implemented and published as two Python packages, Excel2SBOL and SBOL2Excel, under a free and open license and published to pypi for easy pip installation and integration with other projects. Specific benefits of each of the packages are described below.

SBOL2Excel

The SBOL2Excel package is designed in a modular way to allow further expansion. The package’s design includes an initial step from SBOL2 to a dictionary structure. This function could be swapped out with other functions to enable the conversion of other formats such as SBOL3 and GenBank to the Excel Templates.

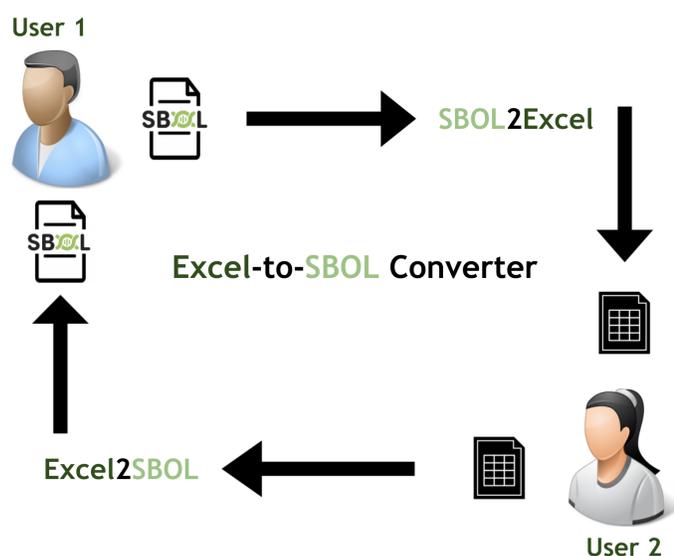


Figure 1: Round trip journey of the Excel-to-SBOL Converter: User 1 converts an SBOL file into an Excel spreadsheet using the SBOL2Excel library, and shares that spreadsheet with User 2, who then converts the Excel spreadsheet back into an SBOL file using the Excel2SBOL library. The SBOL file that User 2 returns to User 1 is exactly the same as the SBOL file that User 1 originally had.

Another key feature is that any arbitrary Component-Definition properties can be converted into Excel columns. However, some common types of properties are given additional post processing to make them more user friendly. SBOL2Excel operates in five general steps:

- (1) An SBOL2 document is read (with property-specific post-processing) and converted to a dictionary.
- (2) The dictionary is converted to a pandas dataframe.
- (3) The dataframe’s columns are reordered based on the default Excel Template.
- (4) Extraneous columns are dropped from the dataframe.
- (5) The dataframe is output into a formatted Excel Spreadsheet.

Figure 2 shows a high level representation of SBOL2Excel’s modular architecture.

Excel2SBOL

This package is able to parse and add information from arbitrary columns into component definitions as annotations. Specific columns are added as a particular type of annotation based on a column conversion table that is accessible in the Excel Template. The package is also able to deal with different templates as long as the Template parameters are added to the modules config file. SBOL2Excel operates in 3 general steps:

- (1) Parse the different sections of the template (Overview Information, Design Description, and Part Table).
- (2) Parse the table indicating how columns are converted.
- (3) For every row, create a component definition and add attributes based on each column in the parts table.

Figure 3 is a high level representation of Excel2SBOL’s modular architecture.

3 DISCUSSION

The next steps in the development of the Excel-SBOL Converter are:

- Making the incorporation of ontologies with Excel Templates easier and less hard-coded, perhaps by relying on RightField Templates. Integrating RightField may be useful in preparing ontologies for templatization.
- Making the benefits of adherence to community standards clearer to users to create greater uptake.
- Establishing a core set of part metadata which can then be enforced via the Excel Templates.
- Making both packages SBOL3 Compliant. Enabling this feature would allow users to perform conversions from SBOL documents (perhaps generated by SynBio-Hub) that contain SBOL2, SBOL3, or both.
- Creating SynBioHub plugins for both packages to integrate their functionality with SynBioHub. This would enable users to perform conversions without directly interacting with Python code.
- Expanding the range of top-level SBOL types that can be converted (e.g. Activities, Composite Components, and Modules).

ACKNOWLEDGEMENTS

JA, JM, and CM are supported by the National Science Foundation under Grant No. 1939892. JM is additionally supported by a Dean’s Graduate Assistantship at the University of Colorado Boulder. IP was supported by the Google Summer of Code. JB is partially supported by AFRL and DARPA contract FA8750-17-C-0184. This document does not contain technology or technical data controlled under either U.S. International Traffic in Arms Regulation or U.S. Export Administration Regulations. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the funding agencies.

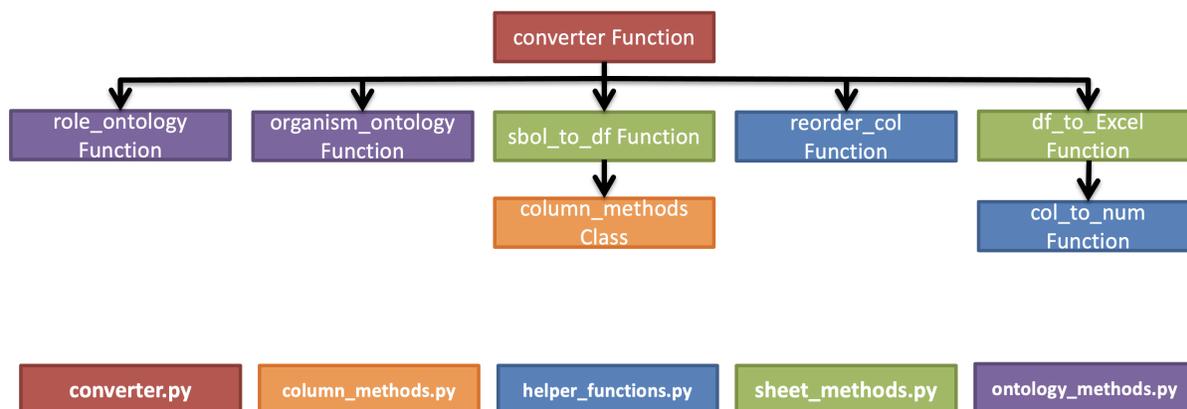


Figure 2: The module architecture of the SBOL2Excel Python Package. Note that the colors of the functions indicate which Python module they can be found in, as indicated in the key on the bottom right.

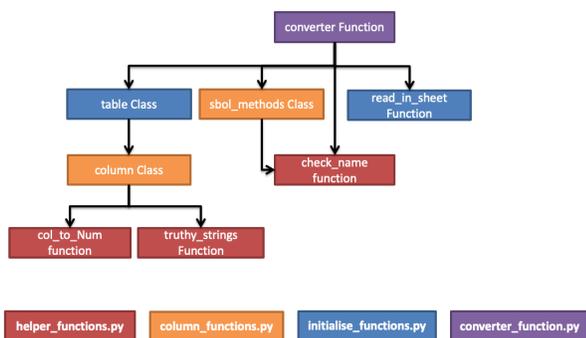


Figure 3: The module architecture of the Excel2SBOL Python Package. Note that the colors of the functions indicate which Python module they can be found in, as indicated in the key on the bottom right.

REFERENCES

- [1] BERGMANN, F. T., ADAMS, R., MOODIE, S., COOPER, J., GLONT, M., GOLEBIEWSKI, M., HUCKA, M., LAIBE, C., MILLER, A. K., NICKERSON, D. P., AND ET AL. Combine archive and omex format: one file to share all information to reproduce a modeling project. *BMC bioinformatics* 15 (Dec 2014), 369.
- [2] BOOTH, I. R. Sysmo: back to the future. *Nature Reviews Microbiology* 5, 8 (Aug 2007), 566–566.
- [3] CANNON, R. C., GLEESON, P., CROOK, S., GANAPATHY, G., MARIN, B., PIASINI, E., AND SILVER, R. A. Lems: a language for expressing complex biological models in concise and hierarchical form and its use in underpinning neuroml 2. *Frontiers in Neuroinformatics* 8 (2014), 79.
- [4] FAEDER, J. R., BLINOV, M. L., AND HLAVACEK, W. S. *Rule-Based Modeling of Biochemical Systems with BioNetGen*, vol. 500 of *Methods in Molecular Biology*. Humana Press, 2009, p. 113–167.
- [5] GALDZICKI, M., CLANCY, K. P., OBERORTNER, E., POCOCK, M., QUINN, J. Y., RODRIGUEZ, C. A., ROEHNER, N., WILSON, M. L., ADAM, L., ANDERSON, J. C., AND ET AL. The synthetic biology open language (sbol) provides a community standard for communicating designs in synthetic biology. *Nature Biotechnology* 32, 6 (Jun 2014), 545–550.
- [6] GENNARI, J. H., NEAL, M. L., GALDZICKI, M., AND COOK, D. L. Multiple ontologies in action: composite annotations for biosimulation models. *Journal of Biomedical Informatics* 44, 1 (Feb 2011), 146–154.
- [7] GLEESON, P., CROOK, S., CANNON, R. C., HINES, M. L., BILLINGS, G. O., FARINELLA, M., MORSE, T. M., DAVISON, A. P., RAY, S., BHALLA, U. S., AND ET AL. Neuroml: A language for describing data driven models of neurons and networks with a high degree of biological detail. *PLOS Computational Biology* 6, 6 (Jun 2010), e1000815.
- [8] HUCKA, M., FINNEY, A., SAURO, H. M., BOLOURI, H., DOYLE, J. C., KITANO, H., ARKIN, A. P., BORNSTEIN, B. J., BRAY, D., CORNISH-BOWDEN, A., AND ET AL. The systems biology markup language (sbml): a medium for representation and exchange of biochemical network models. *Bioinformatics (Oxford, England)* 19, 4 (Mar 2003), 524–531.
- [9] HUCKA, M., NICKERSON, D. P., BADER, G. D., BERGMANN, F. T., COOPER, J., DEMIR, E., GARNY, A., GOLEBIEWSKI, M., MYERS, C. J., SCHREIBER, F., AND ET AL. Promoting coordinated development of community-based information standards for modeling in biology: The combine initiative. *Frontiers in Bioengineering and Biotechnology* 3 (2015), 19.
- [10] KHALIL, A. S., AND COLLINS, J. J. Synthetic biology: applications come of age. *Nature Reviews Genetics* 11, 5 (2010), 367–379.
- [11] MACCAGNAN, A., RIVA, M., FELTRIN, E., SIMIONATI, B., VARDANEGA, T., VALLE, G., AND CANNATA, N. Combining ontologies and workflows to design formal protocols for biological laboratories. *Automated Experimentation* 2, 1 (Apr 2010), 3.
- [12] MCLAUGHLIN, J. A., MYERS, C. J., ZUNDEL, Z., MISIRLI, G., ZHANG, M., OFITERU, I. D., GOÑI-MORENO, A., AND WIPAT, A. Synbiohub: A standards-enabled design repository for synthetic biology. *ACS Synthetic Biology* 7, 2 (Feb 2018), 682–688.
- [13] NICKERSON, D., ATALAG, K., DE BONO, B., GEIGER, J., GOBLE, C., HOLLMANN, S., LONJEN, J., MÜLLER, W., REGIERER, B., STANFORD, N. J., AND ET AL. The human physiome: how standards, software and innovative service infrastructures are providing the building blocks to make it achievable. *Interface Focus* 6, 2 (Apr 2016), 20150103.
- [14] NOVÈRE, N. L., HUCKA, M., MI, H., MOODIE, S., SCHREIBER, F., SOROKIN, A., DEMIR, E., WEGNER, K., ALADJEM, M. I., WIMALARATNE, S. M., AND ET AL. The systems biology graphical notation. *Nature Biotechnology* 27, 8 (Aug 2009), 735–741.
- [15] SPELLMAN, P. T., MILLER, M., STEWART, J., TROUP, C., SARKANS, U., CHERVITZ, S., BERNHART, D., SHERLOCK, G., BALL, C., LEPAGE, M., AND

- ET AL. Design and implementation of microarray gene expression markup language (mage-ml). *Genome Biology* 3, 9 (Aug 2002), research0046.1.
- [16] URQUIZA-GARCÍA, U., ZIELIŃSKI, T., AND MILLAR, A. J. Better research by efficient sharing: evaluation of free management platforms for synthetic biology designs. *Synthetic Biology* 4, 1 (2019), ysz016.
- [17] WALTEMATH, D., ADAMS, R., BERGMANN, F. T., HUCKA, M., KOLPAKOV, F., MILLER, A. K., MORARU, I. I., NICKERSON, D., SAHLE, S., SNOEP, J. L., AND ET AL. Reproducible computational biology experiments with sed-ml—the simulation experiment description markup language. *BMC systems biology* 5 (Dec 2011), 198.
- [18] WITTIG, U., KANIA, R., GOLEBIEWSKI, M., REY, M., SHI, L., JONG, L., ALGAA, E., WEIDEMANN, A., SAUER-DANZWITZ, H., MIR, S., AND ET AL. Sabio-rk—database for biochemical reaction kinetics. *Nucleic Acids Research* 40, Database issue (Jan 2012), D790–796.
- [19] WOLSTENCROFT, K., KREBS, O., SNOEP, J. L., STANFORD, N. J., BACALL, F., GOLEBIEWSKI, M., KUZYAKIV, R., NGUYEN, Q., OWEN, S., SOILAND-REYES, S., AND ET AL. Fairdomhub: a repository and collaboration environment for sharing systems biology research. *Nucleic Acids Research* 45, D1 (Jan 2017), D404–D407.
- [20] WOLSTENCROFT, K., OWEN, S., HORRIDGE, M., JUPP, S., KREBS, O., SNOEP, J., DU PREEZ, F., MUELLER, W., STEVENS, R., AND GOBLE, C. Stealthy annotation of experimental biology by spreadsheets. *Concurrency and Computation: Practice and Experience* 25, 4 (2013), 467–480.
- [21] WOLSTENCROFT, K., OWEN, S., HORRIDGE, M., KREBS, O., MUELLER, W., SNOEP, J. L., DU PREEZ, F., AND GOBLE, C. Rightfield: embedding ontology annotation in spreadsheets. *Bioinformatics* 27, 14 (2011), 2021–2022.
- [22] YU, T., LLOYD, C. M., NICKERSON, D. P., COOLING, M. T., MILLER, A. K., GARNY, A., TERKILDSEN, J. R., LAWSON, J., BRITTEN, R. D., HUNTER, P. J., AND ET AL. The physiome model repository 2. *Bioinformatics (Oxford, England)* 27, 5 (Mar 2011), 743–744.

LOICA: Logical Operators for Integrated Cell Algorithms

Gonzalo Vidal

Pontificia Universidad Católica de Chile
Santiago, Chile
gsvidal@uc.cl

Carlos Vidal-Céspedes

Pontificia Universidad Católica de Chile
Santiago, Chile
carlos.vidal.c@ug.uchile.cl

Timothy James Rudge

Pontificia Universidad Católica de Chile
Santiago, Chile
trudge@uc.cl

1 INTRODUCTION

Synthetic Biology is an interdisciplinary field that mixes life sciences and engineering. From this perspective living systems are objects to engineer, and a rational way to design them is by modifying their genetic code. This can be done by introducing synthetic DNA that encodes a synthetic regulatory network, also known as a genetic circuit. The design-build-test-learn (DBTL) cycle is central to engineering disciplines and each phase requires appropriate tools, standards and workflows, which are still in development in synthetic biology. The synthetic biology open language (SBOL) is an open standard for the representation of *in silico* biological designs that covers the DBTL cycle and has attracted a community of developers that have produced an ecosystem of software tools [4].

Modelling is key to the DBTL cycle and is essential to the design and learn stages; a model states a well-defined hypothesis about the system operation. Abstraction enables the construction and analysis of models based on components, devices, and systems that can be used to compose genetic circuits. It is the basis for genetic design automation (GDA), which can accelerate and automate the genetic circuit design process. In order for GDA to proceed in a rational way, the abstract elements of genetic circuits must be accessible to characterization, allowing parameterization of models of their operation and interactions.

Functional abstraction of DNA sequences as parts such as transcriptional promoters (Pro), ribosome binding sites (RBS), coding sequences (CDS), terminators (Ter) and other elements has enabled the assembly of relatively small genetic circuits [1–3]. However, for large-scale genetic circuit design higher-level abstractions are required, as provided by the logic formalism [6]. In this approach circuit compositions are abstracted into genetic logic gates that transition between discrete low and high steady-state gene expression levels according to input signals, either external or internal to the circuit [9]. These genetic logic circuits can be designed automatically, in an analogous way to electronic circuits, based on the required discrete logical truth table [6], however this specification requires knowledge of the domain-specific programming language *Verilog*.

Despite the discrete logical design formalism, these genetic circuits are dynamical systems and can have autonomous, continuous non-steady-state dynamics, displaying complex and rich behaviors from bi-stability to oscillations and even chaos [2, 3, 11]. Furthermore, typical operating conditions for engineered circuits like colonies, bioreactors or gut microbiomes are time varying, which can lead to complex behaviors from even simple genetic circuits [7].

To design genetic circuit temporal dynamics we therefore require kinetic gene expression data generated at the test phase. This data must be integrated with models to enable characterization of abstracted parts, devices and systems, as well as metadata, including the DNA part composition, to enable automated design. Thus there is a need for software design tools that integrate abstract circuit designs, dynamical models, kinetic gene expression data, and DNA part composition via common exchange standards in a user-friendly and accessible fashion.

2 RESULTS

Logical Operators for Integrated Cell Algorithms (LOICA) provides a high-level genetic design abstraction using a simple and flexible object-oriented programming approach in Python. LOICA integrates models with experimental data via two-way communication with Flapjack, a data management and analysis tool for genetic circuit characterization [12]. This communication not only provides direct access to experimental characterization data, but also to DNA design composition and sequence via SBOL contained in SynBioHub [5], enabling characterization and simulation in the same tool, but also facilitating exchange with other tools such as iBioSim [11].

The basic objects in LOICA are *Operator* and *GeneProduct*, which may be either a *Regulator* or *Reporter* (Figure 1A). A *Regulator* represents a molecular species that regulates gene expression. A *Reporter* is a molecular species that provides a measurable signal, such as a fluorescent protein. The *Operator* maps one or more *Regulator* concentrations to one or more *GeneProduct* synthesis rates. An *Operator* can be implemented in DNA as a combination of Pro and RBS, and the *Regulator* could be a CDS of a transcription factor or of a

regulatory RNA. The interactions between the *Operators* and the *Regulators* encode models for genetic circuit temporal dynamics, which are simulated with differential equations. The system is thus:

$$\frac{d\mathbf{p}}{dt} = \Psi(\mathbf{r}) - \Gamma\mathbf{p} - \mu(t)\mathbf{p}, \quad (1)$$

$$\Psi(\mathbf{r}) = \sum_k \Phi_k(\mathbf{r}), \quad (2)$$

where $\mathbf{p} = (p_0, p_1, \dots, p_{N-1})^T$ is the vector of *GeneProducts*, which includes different *Regulators* ($\mathbf{r} = (r_0, r_1, \dots, r_{M-1})^T$) and *Reporters* ($\mathbf{s} = (s_0, s_1, \dots, s_{N-M-1})^T$). The non-linear operator Ψ maps *Regulator* concentrations to *GeneProduct* synthesis rates. Γ is a diagonal matrix of *GeneProduct* degradation rates, and $\mu(t)$ is the instantaneous growth rate of the cells. Equation 1 shows the overall system where Ψ encodes the whole circuit, and consists of a sum of individual LOICA *Operators* Φ_k (Equation 2).

Figure 1B shows a mathematical model that results from the interaction of an *Operator* encoding a simple NOT logic modeled with a Hill equation as transfer function. It can be implemented as a promoter containing repressor binding sites combined with a ribosome binding site (RBS). The logical *Operator* can thus be instantiated as a genetic device that is repressed by an input *Regulator* and outputs a *GeneProduct* synthesis rate. Note that LOICA can be used to define an *Operator* as any operation that maps from input *Regulator* concentrations to output synthesis rates.

The repressilator is a useful dynamical system case study because it produces continuous sustained oscillations that escapes ON/OFF logic [2]. To model it, we consider a simple balanced ring oscillator with three NOT *Operators* connected with three different *Regulators*. The *Operators* and *Regulators* are incorporated into a *GeneticNetwork*, linked with a Flapjack *Vector*, which with the *Metabolism* drives the dynamics of the *Sample*, also corresponding to a Flapjack *Sample* (Figure 1A). For the circuit to produce a measurable signal we add three *Operators* using the same inputs but changing the outputs to three different *Reporters*, linked with the Flapjack *Signal* model (Figure 2A). The code to generate this model is in Figure 2B. This approach is used to generate synthetic data (a LOICA *Assay*) from models that can be uploaded to Flapjack. It is then easy to access Flapjack’s genetic circuit characterization tools, data management and data visualization through a Python package (pyFlapjack) or the web interface shown in Figure 2D.

We have described how to use LOICA to generate and analyze simulation data from models. Another example workflow goes from data to model parameterization and could be as follows for the learn stage of the DBTL cycle. First, a *GeneticNetwork* is assembled from a collection of *Operators*, linked by various *Regulators* and *Reporters* in some topology.

Each *Operator* is then linked to experimental data contained in Flapjack, which corresponds to measurements of the auxiliary circuits required to parameterize the model encoded. For example, a NOT *Operator* links to data of a chemical signal receiver and a chemical signal inverter measured in a range of signal concentrations (LOICA and Flapjack *Supplement*). The *Operator* provides a function that then extracts the data from Flapjack and uses it to parameterize the corresponding model (Figure 1B). In the example shown that means fitting parameters a, b, K, n by least squares minimization of the difference between the experimental data and the solution of differential equation models of the auxiliary circuits (Equation 1). Each *Operator* thus contains the information required to characterize itself. Therefore with LOICA data driven models that include mathematical constraints it is possible to design circuit topologies and look for likely functional designs initiating a new round of the DBTL cycle.

3 CONCLUSION

LOICA integrates the design and characterization of genetic circuit dynamics into Python workspaces, providing an easy-to-understand design abstraction implemented using simple object-oriented programming principles. This programming interface does not require specialist or domain-specific knowledge, but leverages common programming skills, making it accessible but also providing customization capabilities for advanced users. LOICA is able to simulate abstract genetic circuit designs using differential equation models. It abstracts genetic circuit designs into objects which are capable of characterizing themselves via links to data in Flapjack. As data relating to genetic components is updated in Flapjack, the fitted parameters can be automatically updated upon characterization. Flapjack also provides the connection to SynBioHub which allows design and characterization based on SBOL.

4 FUTURE WORK

We aim to complete and automate the DBTL cycle in synthetic biology, proposing a workflow that integrates LOICA, Flapjack [12] and SynBioHub [5] via SBOL3 [4]. In the design stage SBOL/SBML will provide flexibility to use other existing bioCAD tools such as iBioSim [11] as part of the workflow. To make a direct transition to the build stage LOICA will generate a human/machine readable protocols for assembly using Opentrons liquid handling robots [10] and their open Python API. We will expand LOICA’s capabilities to single-cell resolution spatio-temporal systems by connecting it to CellModeller [8] for individual based modelling and stochastic simulations. For the test stage we will develop open science hardware for measuring genetic circuit dynamics that could be used to obtain kinetic data for direct upload to Flapjack.

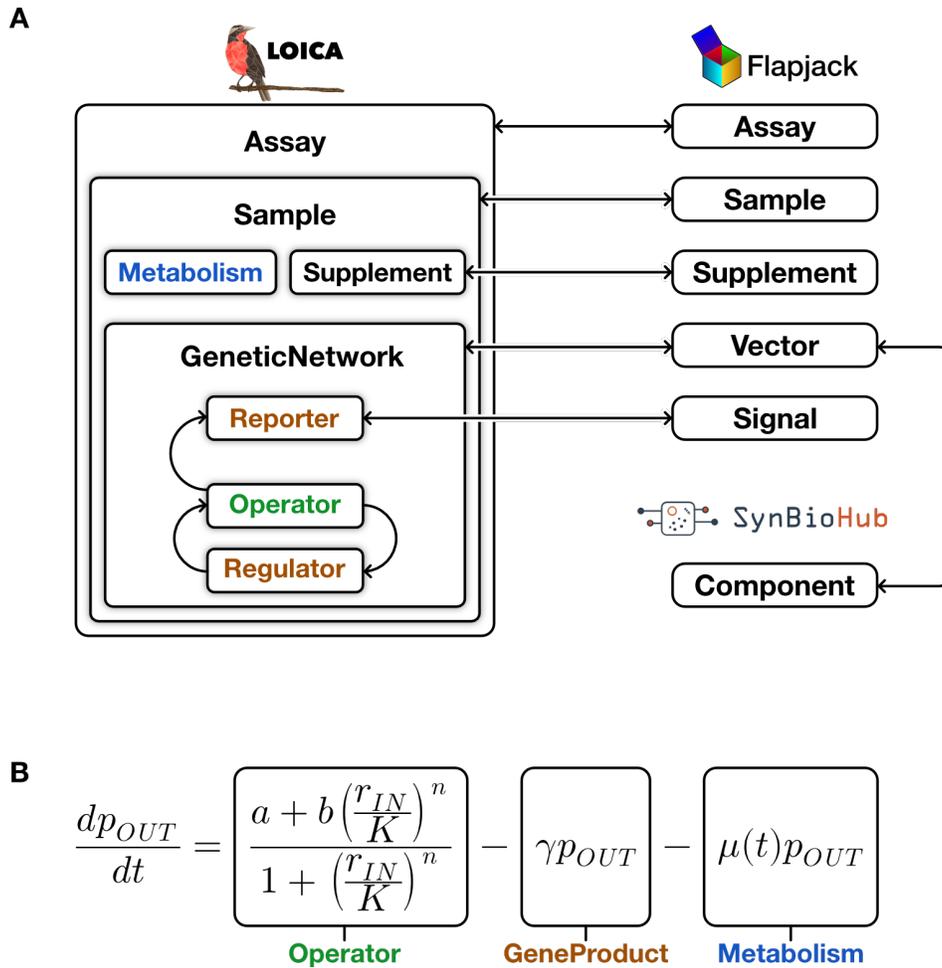


Figure 1: Diagram of model generation in LOICA. A. Diagram of an Assay encapsulating a Sample which in turn encapsulates Metabolism, Supplement, and GeneticNetwork. In the later, the Operator and Regulator are interacting to generate a model. On the right side the different interactions with the Flapjack and SynBioHub models are shown. B. Mathematical model of a NOT Operator with one input and one output generated through LOICA object interactions. Here p_{OUT} is a GeneProduct (Regulator or Reporter), output of the Operator. In the Operator a is the basal or leaky gene expression, b is the regulated gene expression, r_{IN} is a Regulator concentration, K is the switching concentration, and n is the cooperativity degree of r_{IN} with respect to the Operator. In the GeneProduct γ is the degradation rate of p_{OUT} . In Metabolism $\mu(t)$ is the instantaneous growth rate which dilutes p_{OUT} . Here the Operator is encoded by a Hill equation transfer function that states its regulation by r_{IN} . The transfer function could be any mapping from input concentration to synthetis rate, making LOICA Operators flexible and easy to extend.

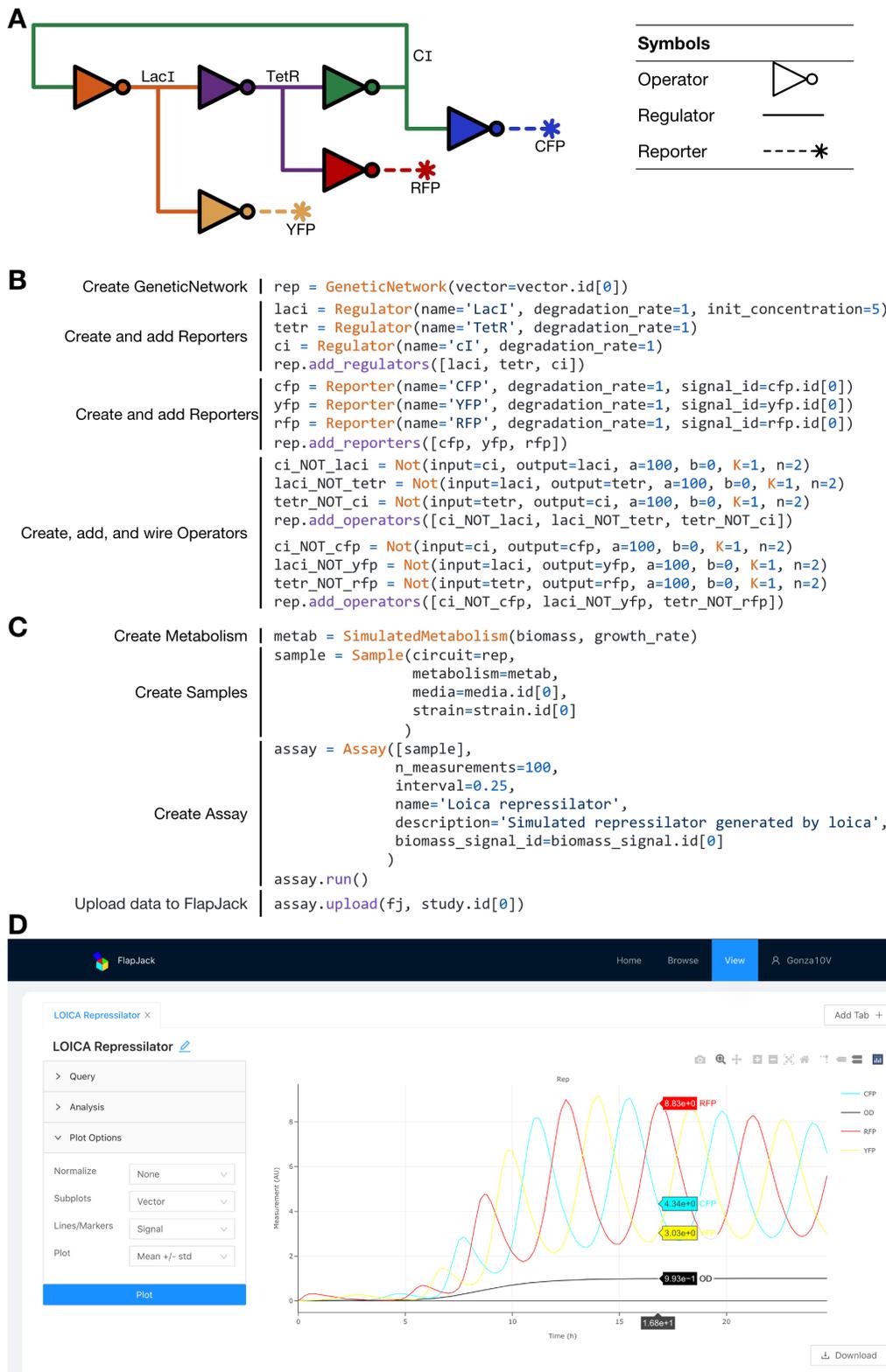


Figure 2: Example of repressilator design in LOICA. A. LOICA diagram of the modeled repressilator circuit and respective symbols. **B.** Python code that generates a repressilator in LOICA. GeneticNetwork construction is the first step where the user states all the objects and their interaction. **C.** Next during Assay setup the user initializes and runs the simulation, and the results can be uploaded to Flapjack. The two-way communication with Flapjack allows data storage and management, enables various analyses to be performed, and allows Operators to characterize themselves. **D.** Data exploration in Flapjack via the web interface. This interface allows to query, analyze and plot the data, which can also be done through the pyFlapjack Python package.

REFERENCES

- [1] DANINO, T., MONDRAGÓN-PALOMINO, O., TSIMRING, L., AND HASTY, J. A synchronized quorum of genetic clocks. *Nature* 463, 7279 (2010), 326–330.
- [2] ELOWITZ, M. B., AND LEIBLER, S. A synthetic oscillatory network of transcriptional regulators. *Nature* 403, 6767 (2000), 335–338.
- [3] GARDNER, T. S., CANTOR, C. R., AND COLLINS, J. J. Construction of a genetic toggle switch in *escherichia coli*. *Nature* 403, 6767 (2000), 339–342.
- [4] McLAUGHLIN, J. A., BEAL, J., MISIRLI, G., GRÜNBERG, R., BARTLEY, B. A., SCOTT-BROWN, J., VAIDYANATHAN, P., FONTANARROSA, P., OBERORTNER, E., WIPAT, A., ET AL. The synthetic biology open language (sbol) version 3: simplified data exchange for bioengineering. *Frontiers in Bioengineering and Biotechnology* 8 (2020), 1009.
- [5] McLAUGHLIN, J. A., MYERS, C. J., ZUNDEL, Z., MISIRLI, G., ZHANG, M., OFITERU, I. D., GONI-MORENO, A., AND WIPAT, A. Synbiohub: a standards-enabled design repository for synthetic biology. *ACS synthetic biology* 7, 2 (2018), 682–688.
- [6] NIELSEN, A. A., DER, B. S., SHIN, J., VAIDYANATHAN, P., PARALANOV, V., STRYCHALSKI, E. A., ROSS, D., DENSMORE, D., AND VOIGT, C. A. Genetic circuit design automation. *Science* 352, 6281 (2016).
- [7] PEÑA, G. A. V., VIDAL-CÉSPEDES, C. I., ET AL. Accurate reconstruction of dynamic gene expression and growth rate profiles from noisy measurements. *bioRxiv* (2021).
- [8] RUDGE, T. J., STEINER, P. J., PHILLIPS, A., AND HASELOFF, J. Computational modeling of synthetic microbial biofilms. *ACS synthetic biology* 1, 8 (2012), 345–352.
- [9] SHIN, J., ZHANG, S., DER, B. S., NIELSEN, A. A., AND VOIGT, C. A. Programming *escherichia coli* to function as a digital display. *Molecular systems biology* 16, 3 (2020), e9401.
- [10] STORCH, M., HAINES, M. C., AND BALDWIN, G. S. Dna-bot: a low-cost, automated dna assembly platform for synthetic biology. *Synthetic Biology* 5, 1 (2020), ysaa010.
- [11] WATANABE, L., NGUYEN, T., ZHANG, M., ZUNDEL, Z., ZHANG, Z., MADSEN, C., ROEHNER, N., AND MYERS, C. ibiosim 3: A tool for model-based genetic circuit design. *ACS Synthetic Biology* 8, 7 (2019), 1560–1563.
- [12] YÁÑEZ FELIÚ, G., EARLE GÓMEZ, B., CODOCEO BERROCAL, V., MUÑOZ SILVA, M., NUÑEZ, I. N., MATUTE, T. F., ARCE MEDINA, A., VIDAL, G., VIDAL CÉSPEDES, C., DAHLIN, J., ET AL. Flapjack: Data management and analysis for genetic circuit characterization. *ACS Synthetic Biology* 10, 1 (2020), 183–191.

Classifying Literature Extracted Events for Automated Model Extension

Casey Hansen¹, Julia Kisslinger², Neal Krishna³, Emilee Holtzapple⁴, Yasmine Ahmed², Natasa Miskov-Zivanov^{1,2,4}

¹University of Pittsburgh, Department of Bioengineering, ²University of Pittsburgh, Department of Electrical and Computer Engineering, ³ University of Connecticut, ⁴ University of Pittsburgh, Department of Computational and Systems Biology

{ceh92,juk77,erh87,yaa38,nmzivanov}@pitt.edu,neal.krishna@uconn.edu

1 INTRODUCTION

Machine reading tools are able to quickly and automatically curate vast amounts of information from relevant published literature [6][2]. This curated information can be used to build biological computational models or expand upon existing models. However, the information gleaned by machine readers is both vast and varied in quality. Machine readers must work to extract standardized biological interactions from inconsistent terminology and complex sentence structures, which sometimes leads to extraction errors.

Previously we have developed VIOLIN (Verifying Interactions of Likely Importance to the Network) a tool to automatically classify and judge biological interactions extracted from relevant literature. With VIOLIN, we are able to take these literature extracted events (LEEs) and compare them to an existing biological model, determining whether a given LEE agrees with the model (corroborates), introduces new information to the model (extends), disputes the model (contradicts), or requires manual review (flagged). Each LEE is assigned four numerical values to represent its relationship to the model system (Match Score), its classification category (Kind Score), its frequency (Evidence Score), and extraction confidence (Epistemic Value). These values are combined into a Total Score to allow for automatic filtering and classification of large sets of LEEs curated from multiple sources. To further increase the utility of VIOLIN, we now seek to integrate VIOLIN as part of an automated model-building framework (Figure 1).

Current approaches towards building and extending models have two major pitfalls. They either focus on only a single step of the process [8][7], or the decision metrics lack depth, focusing on the machine reading output or model as separate entities, more than the relationship between the two [2].

We first integrated VIOLIN with the filtering tool FLUTE (FiLter for Understanding True Events) [3], to make use of the expert data gathered in public databases. The FLUTE tool connects to public protein interaction databases to judge the accuracy of an LEE. This integration allows us to balance the removal of erroneous extractions while retaining novel interactions which may not yet be represented in a database.

We next integrated VIOLIN with CLARINET (CLARifying NETWORKS) [1], an automated model extension tool. Where VIOLIN classifies individual LEEs for their relevance and usefulness to a given model, CLARINET classifies candidate extensions as clusters of biological interactions, taking into account how LEEs are connected to each other, in addition to their connection to the baseline model.

These three tools together create a powerful method of taking information-rich relevant literature and identifying the highest quality events for extending a baseline model.

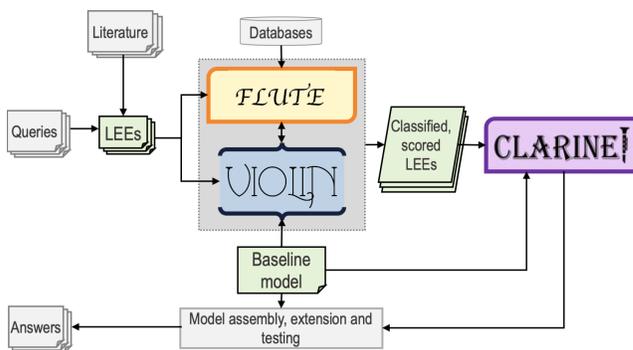


Figure 1: An outline of automated information extraction and the model assembly, highlighting the roles FLUTE, VIOLIN, and CLARINET hold: FLUTE and VIOLIN judge the quality, relevance, and usefulness of LEEs on an individual basis, and then CLARINET judges how the LEEs connect, both to each other and the baseline model.

2 METHODS

To evaluate the integration of VIOLIN and FLUTE, we used the following as inputs (1) three computational models, namely a model of Skel-133 Melanoma, a model of human T-cells [4], and a model of the BDNF pathway as it relates to Major Depressive Disorder (MDD) [5], and (2) four LEE sets for each model. From these inputs, we generated three types of outputs: LEE sets classified by VIOLIN only (control), LEE sets first filtered by FLUTE and then classified by VIOLIN (pre-processed), and LEE sets first classified by VIOLIN and then filtered by FLUTE (post-processed).

We next investigated the integration of VIOLIN and CLARINET, using a Glioblastoma Multiforme model and two highly specialized LEE sets. Our first LEE set (R_{G1}) contained 10,130 LEEs from 242 papers, and the second (R_{G2}) contained 25,875 LEEs from 454 papers. From these inputs, we also created three data outputs: candidate clusters created from the raw LEE sets (control), candidate clusters created from the Total VIOLIN output, which lists only the unique LEEs (unique), and candidate clusters created from only the VIOLIN extensions (extensions). Table 1 shows a summary of the input parameters for both parts of our investigation.

Table 1: Testing Inputs

LEE Suffix	Model	Model Nodes	LEE sets
R _A	T cell	61	4
R _B	Melanoma	225	4
R _C	MDD	72	4
R _G	GBM	238	2

3 RESULTS

For our VIOLIN-FLUTE integration, we found that post-processing methods consistently retained the greatest number of LEEs, and that this is true across all VIOLIN classification categories (Figure 2). Contradictions and Extensions had the lowest average retention rate, and flagged had the highest (Table 2). This supports our previous suggestion that the contradiction category can be used to filter out machine reading errors. As expected, those LEEs with high evidence scores are retained more often than those with low evidence scores.

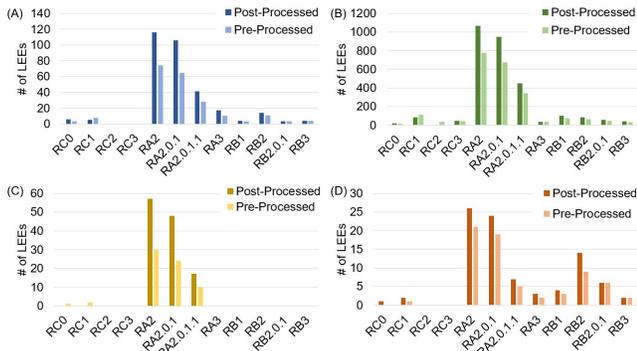


Figure 2: Retention counts for each VIOLIN classification: (A) shows corroborations, (B) shows extensions, (C) shows contradictions, and (D) shows flagged

Table 2: Average Retention Rates

	Pre-Processed %	Post-Processed %
Corroborations	32.9	42.1
Extensions	21.7	26.9
Contradictions	10.7	16.6
Flagged	39.3	61.6

For the VIOLIN-CLARINET integration, we observed the size and central nodes of the candidate clusters for the control, unique, and extension output from CLARINET (Figure 3). We found that just the act of comparing the control output to the unique output, which contains only single instances of a given LEE, has an effect on the outcome of the candidate clusters. The candidate clusters from the extensions are an even more focused input, as they only present LEEs for consideration which are known to present new information to the model. This suggests that forming candidate

clusters from raw machine reading output is influenced by corroborative or contradictory LEEs, as well as machine reading errors, and having more directed LEE sets would produce more directed clusters.

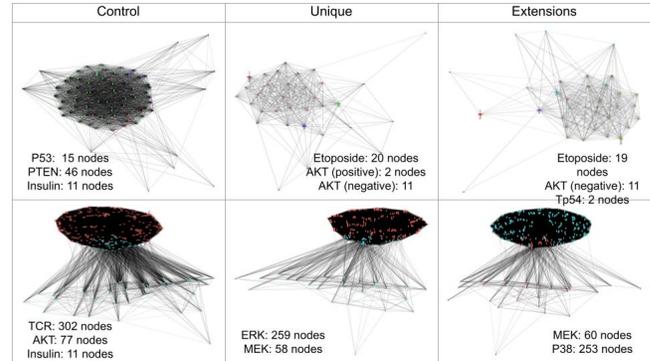


Figure 3: Candidate clusters for the control, unique, and extensions input compared to the GBM model using CLARINET. The top row was created from the R_{G1} LEE set, and the bottom row was created from the R_{G2} LEE set

4 CONCLUSIONS

Integrating VIOLIN with FLUTE and CLARINET showed us the promising outcome of combining these individually effective tools. The options of using FLUTE with VIOLIN allows the user to determine the importance of removing erroneous information versus retaining novel interactions. Our results from CLARINET show that narrowing an LEE set down to those which are most useful for extension changes the candidate extensions, and VIOLIN allows this process to be fast and automatic. Our next step is to further investigate approaches to utilize the integration of these three tools towards automated creation of useful and reliable models.

ACKNOWLEDGEMENTS

This project was funded by DARPA award W911NF-17-1-0135.

REFERENCES

- [1] AHMED, Y., TELMER, C., AND MISKOV-ZIVANOV, N. CLARINET: Efficient learning of dynamic network models from literature.
- [2] GYORI, B. M., BACHMAN, J. A., SUBRAMANIAN, K., MUHLICH, J. L., GALESCU, L., AND SORGER, P. K. From word models to executable models of signaling networks using automated assembly. 954.
- [3] HOLTZAPPLE, E., TELMER, C. A., AND MISKOV-ZIVANOV, N. FLUTE: Fast and reliable knowledge retrieval from biomedical literature. [_eprint: https://academic.oup.com/database/article-pdf/doi/10.1093/database/baaa056/33571972/baaa056.pdf](https://academic.oup.com/database/article-pdf/doi/10.1093/database/baaa056/33571972/baaa056.pdf).
- [4] MISKOV-ZIVANOV, N., TURNER, M. S., KANE, L. P., MOREL, P. A., AND FAEDER, J. R. The duration of t cell stimulation is a critical determinant of cell fate and plasticity. *ra97-ra97*.
- [5] SANDHYA, V. K., RAJU, R., VERMA, R., ADVANI, J., SHARMA, R., RADHAKRISHNAN, A., NANJAPPA, V., NARAYANA, J., SOMANI, B. L., MUKHERJEE, K. K., PANDEY, A., CHRISTOPHER, R., AND PRASAD, T. S. K. A network map of BDNF/TRKB and BDNF/p75ntr signaling system. 301–307.
- [6] VALENZUELA-ESCARCEGA, M. A., HAHN-POWELL, G., SURDEANU, M., AND HICKS, T. A domain-independent rule-based framework for event extraction. In *Proceedings of ACL-IJCNLP 2015 System Demonstrations* (Beijing, China, July 2015), Association for Computational Linguistics and The Asian Federation of Natural Language Processing, pp. 127–132.

- [7] VON MERING, C., JENSEN, L. J., SNEL, B., HOOPER, S. D., KRUPP, M., FOGHERINI, M., JOUFFRE, N., HUYNEN, M. A., AND BORK, P. STRING: known and predicted protein-protein associations, integrated and transferred across organisms. *Nucleic Acids Res* 33, Database issue (Jan 2005), D433–437.
- [8] ZERVA, C., BATISTA-NAVARRO, R., DAY, P., AND ANANIADOU, S. Using uncertainty to link and rank evidence from biomedical literature for model curation. *Bioinformatics* 33, 23 (Dec 2017), 3784–3792.

New advances in the automation of context-aware information selection and guided model assembly

Yasmine Ahmed¹, Adam A. Butchy², Khaled Sayed¹, Cheryl Telmer³, Natasa Miskov-Zivanov^{1,2,4}

¹Department of Electrical and Computer Engineering, ²Department of Bioengineering, ⁴Department of Computational and Systems Biology, University of Pittsburgh, ³Molecular Biosensor and Imaging Center, Carnegie Mellon University {YAA38,Adam.Butchy,KSS60,NMZivanov}@pitt.edu,Ctelmer@cmu.edu

1 INTRODUCTION

Modeling complex systems or extending existing models with new information enables a better understanding of these systems [5]. New information can be extracted from different knowledge sources—such as expert knowledge, published literature and pathway databases—and used to assemble or extend models (Figure 1 (a)). However, modeling is a time and labor-intensive task, often limited by the knowledge and experience of the modelers. With new research articles published each day, there is a pressing need for an automated method that updates models with new information efficiently and automatically, while preserving the usefulness and accuracy of the original models. Recently, there has been a push in the field of synthetic biology to automate the entire pathway of model assembly, starting with collecting biological interactions, assembling a model, and performing simulations [3]. A typical model assembly pipeline (Figure 1 (b)) begins with a question about the system under study. This question is converted into a search engine query to identify and extract the most relevant papers. Biological events are extracted from those papers and used to assemble or extend models. The newly assembled models are then analyzed and evaluated to determine if they satisfy desired system behavior. In this work, we survey the most recent automated model assembly efforts. Specifically, we will review five tools: Layer-based [8], Genetic Algorithm (GA) based [12], ACCORDION [1], CLARINET [2] and FIDDLE [4]. We will emphasize the applicability and benefits of each tool using a case study of T cell differentiation model [6] [9].

2 BACKGROUND

Cellular signaling pathways can be modeled as directed graphs, with nodes representing pathway elements, and edges representing interactions between elements. To study the dynamics of such systems, all the presented tools use executable models, where discrete variables represent states of model elements, and each element can have a state transition function or update function. A baseline model, and an output from a machine reading engine are the inputs to the model assembly pipeline. Each model assembly pipeline generates candidate models of the system under study. Model checking is then used to verify whether each candidate model

satisfies a set of properties describing expected behavior of the system. Here, we compare the tools using the same T cell model described in [6], and suggested set of interactions from an open-source reading engine, REACH [15]. We expressed both the model and the reading output using an element-based BioRECIPES format [14] and we used the DiSH simulator [13] to observe dynamic behavior of the baseline and newly assembled models. We also used statistical model checking [17] [7] [10] [11] to test all generated models against formally defined properties.

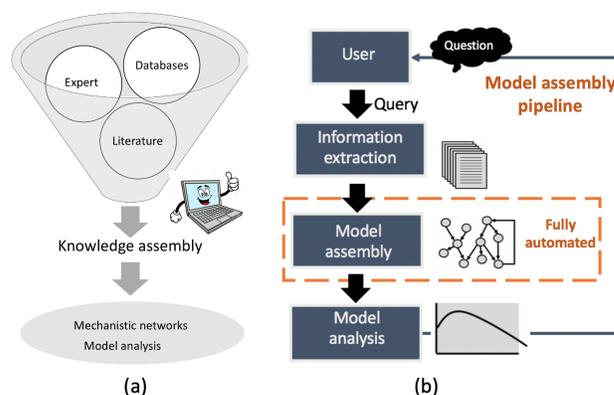


Figure 1: (a) Knowledge assembly, conceptual overview; (b) Automated model assembly pipeline.

3 AUTOMATED MODEL ASSEMBLY TOOLS

3.1 Layer-based approach. In [8], the authors proposed a method that starts with a baseline model and selects interactions extracted from published literature automatically. The proposed method groups the information extracted from literature into layers, based on their proximity to the important elements in the baseline model. The pieces of information organized in such layers are then added to the baseline model, so that the extended model satisfies predefined system properties. The proposed method helps identify some new interactions without trying the extracted interactions all at once or one interaction at a time. Since the extension method adds new interactions based on their proximity to existing models, this method becomes impractical with large-scale models.

3.2 GA-based approach. Another model extension method that uses a Genetic Algorithm [18] was proposed in [12]. The authors in [12] removed a group of elements from an existing model in a random way and they mixed them with randomly created interactions to mimic the output of machine reading engines. Eventually, they applied the genetic algorithm to search for the extensions that optimally reconstructed the model. It has been proved in [12] that the GA-based approach was able to extract a set of extensions that led to the desired system behavior. The main disadvantages of the GA-based approach include non-determinism, as the solution may vary across multiple algorithm executions on the same inputs, as well as issues with scalability.

3.3 ACCORDION (Automated Clustering Conditional On Relating Data of Interactions tO a Network). A tool that automatically and efficiently assembles the information extracted from available literature into models, evaluates the dynamic behavior of newly assembled models, and selects the most suitable model to address user questions as described in [1]. In contrast to [8] and [12], ACCORDION focuses on identifying clusters of strongly connected elements in the newly extracted information that have a measurable impact when added to the model. ACCORDION uses Markov Clustering algorithm (MCL) [16], an unsupervised graph clustering algorithm, to group interactions obtained from literature by machine reading. Eventually, it finds return paths that start in the baseline model, go through one or more clusters, and end in the baseline model; the baseline model and the clusters on such return path form a candidate model.

3.4 CLARINET (CLARIfying NETworks). Recently, a novel methodology was proposed in [2] to expand dynamic network models using the information extracted from published literature by machine reading engines. CLARINET organizes the extracted events as a collaboration graph and uses several novel metrics for evaluating these events individually, in pairs, and in groups. The metrics introduced by CLARINET are based on the frequency of occurrence and co-occurrence of events in literature, and their connectivity to the baseline model. CLARINET is scalable; its average runtime is at the order of seconds when processing several thousand interactions.

3.5 FIDDLE (Finding Interactions using Diagram Driven modeL Extension). A tool described in [4] that employs two methods based on network search algorithms—Breadth First Addition (BFA) and Depth First Addition (DFA)—to automatically assemble or extend models with the knowledge extracted from published literature. FIDDLE is able to refine

models by systematically adding known biological interactions into intermediate models, measuring changes in model performance, and then adding or discarding interactions based on whether they improve the model performance metric. Both BFA and DFA scale linearly with the size of the model they are tasked to extend, and the number of potential interactions with which to extend the model.

4 RESULTS AND DISCUSSION

To demonstrate the accuracy, efficiency, and utility of each tool, we have selected a computational model of T cell differentiation [9]. Our main goal with this case study is to show that each tool is able to automatically assemble and extend an existing published model into another published and manually built model using new elements and new interactions automatically extracted from published literature. As the final golden model, we used the T cell model published in [6] and the set of desired system properties discussed in [9] and [6]. The complete list of 27 properties is shown in Table 1. The golden model and the properties are used to evaluate the automatically assembled model obtained by each tool. Figure 2 highlights the differences between the results obtained for each tool when tested using statistical model checking. The GA-based method features the best performance as scored through statistical model checking. Due to its iterative nature, the time required to perform GA-based extension increases with the number of possible extensions, and can be prohibitively long when applied to large scale models [2]. Both ACCORDION and CLARINET balance performance with scalability and can be applied to large scale models, as well as large scale machine reading output, as demonstrated in [1] and [2]. CLARINET scores the newly extracted events based on both the evidence from literature and the connectivity to the baseline model. If the user is interested in collecting new interactions that are strongly connected to each other and strongly connected to the baseline model, then, ACCORDION would be a better choice; since it adds paths of connected interactions, which are at the same time connected to the baseline model. The layer-based and BFA methods perform similarly, despite adding different number of extensions to the baseline model. The layer-based method is meant to be applied when the user is interested in collecting new, relevant interactions that are directly connected to the baseline model. The DFA method performed the worst, scoring below the baseline model. This can be attributed to optimizing a scoring metric different than statistical model checking. In fact, both FIDDLE methods attempt to optimize the same metric with the fewest number of extensions to the baseline model. Their poor performance points to their metric being a poor stand in for statistical model checking, and the stipulation to minimize the number of additional extensions as an unnecessary restraint.

Table 1: Set of properties that are observed to be true in T cells [6] [9].

Prop#	Description
Scenario 0: No TCR	
1	Once deactivated, AKT will remain inactive until end of analyzed period
2	Once activated, PTEN will remain active until end of analyzed period
3	Once deactivated, FOXP3 will remain inactive until end of analyzed period
4	Once deactivated, IL2 will remain inactive until end of analyzed period
5	Once deactivated, CD25 will remain inactive until end of analyzed period
6	Once deactivated, STAT5 will remain inactive until end of analyzed period
7	Once deactivated, mTOR will remain inactive until end of analyzed period
8	Once deactivated, mTORC2 will remain inactive until end of analyzed period
9	Once activated, FOXO1 will remain active until end of analyzed period
Scenario 1: Low TCR	
10	Once deactivated, AKT will remain inactive until end of analyzed period
11	Once activated, PTEN will remain active until end of analyzed period
12	Once activated, FOXP3 will remain active until end of analyzed period
13	Once deactivated, IL2 will remain inactive until end of analyzed period
14	Once activated, CD25 will remain active until end of analyzed period
15	Once activated, STAT5 will remain active until end of analyzed period
16	Once activated, mTOR will remain active until end of analyzed period
17	Once activated, mTORC2 will remain active until end of analyzed period
18	Once activated, FOXO1 will remain active until end of analyzed period
Scenario 2: High TCR	
19	Once deactivated, AKT will remain inactive until end of analyzed period
20	In developing Th, PTEN decreases and remains absent
21	Once deactivated, FOXP3 will remain inactive until end of analyzed period
22	Once activated, IL2 will remain active until end of analyzed period
23	Once activated, CD25 will remain active until end of analyzed period
24	Once activated, STAT5 will remain active until end of analyzed period
25	Once deactivated, mTOR will remain inactive until end of analyzed period
26	Once activated, mTORC2 will remain active until end of analyzed period
27	Once activated, FOXO1 will remain active until end of analyzed period

5 CONCLUSION AND FUTURE WORK

Automatically extending models with the information published in literature allows for rapid collection of the existing information in a consistent and comprehensive way. It also facilitates information reuse and data reproducibility. In this review, we described five recent efforts in this direction. We demonstrated the respective benefits and drawbacks of each tool and we tested them on a previously published biological model. These methods and software tools represent a novel effort to replace hundreds or thousands of manual experiments, and have a potential to significantly accelerate the advancement of scientific knowledge. As our future work, we will conduct a more in-depth comparison of the five tools to even more precisely evaluate their advantages and drawbacks. We plan to apply the proposed methods on several other models in different biological domain, and we will work on parallelization of the model checking algorithm to further increase its execution efficiency.

6 ACKNOWLEDGEMENT

This project was funded by DARPA award W911NF-17-1-0135.

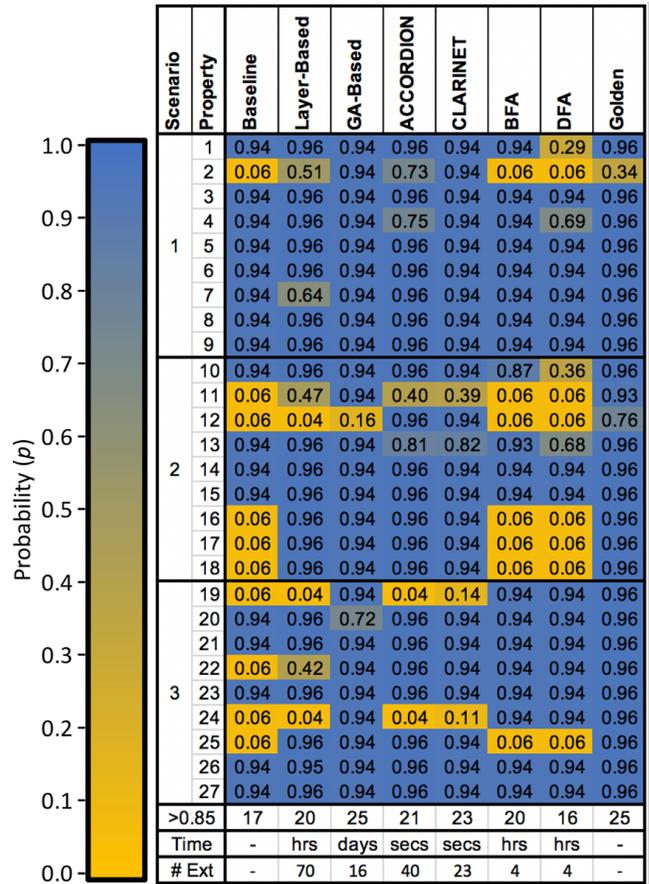


Figure 2: Comparison of the model checking probability estimates p for the baseline model, golden model, and the best model obtained from each of the five tools: Layer-based, GA-based, ACCORDION, CLARINET and FIDDLE (BFA and DFA), when run on a 3.3 GHz Intel Core i5 processor. In the last three rows, we show the number of properties with probability estimates >0.85 , the length of time for each method, and the number of extensions added to the baseline model.

REFERENCES

- [1] AHMED, Y., TELMER, C., AND MISKOV-ZIVANOV, N. Accordion: Clustering and selecting relevant data for guided network extension and query answering, 2020.
- [2] AHMED, Y., TELMER, C., AND MISKOV-ZIVANOV, N. CLARINET: Efficient learning of dynamic network models from literature. *Bioinformatics Advances* (June 2021).
- [3] APPLETON, E., MADSEN, C., ROEHNER, N., AND DENSMORE, D. Design automation in synthetic biology. *Pubmed*.
- [4] BUTCHY, A. A., TELMER, C. A., AND MISKOV-ZIVANOV, N. FIDDLE: Efficient assembly of networks that satisfy desired behavior.
- [5] FISHER, J., AND HENZINGER, T. A. Executable cell biology. *Nature Biotechnology* 25, 11 (Nov. 2007), 1239–1249.
- [6] HAWSE, W. F., SHEEHAN, R. P., MISKOV-ZIVANOV, N., MENK, A. V., KANE, L. P., FAEDER, J. R., AND MOREL, P. A. Cutting edge: Differential regulation of PTEN by TCR, akt, and FoxO1 controls CD4 t cell fate

- decisions. *The Journal of Immunology* 194, 10 (Apr. 2015), 4615–4619.
- [7] JHA, S. K., CLARKE, E. M., LANGMEAD, C. J., LEGAY, A., PLATZER, A., AND ZULIANI, P. A bayesian approach to model checking biological systems. In *Computational Methods in Systems Biology*. Springer Berlin Heidelberg, 2009, pp. 218–234.
- [8] LIANG, K.-W., WANG, Q., TELMER, C., RAVICHANDRAN, D., SPIRITES, P., AND MISKOV-ZIVANOV, N. Methods to expand cell signaling models using automated reading and model checking. In *Computational Methods in Systems Biology*. Springer International Publishing, 2017, pp. 145–159.
- [9] MISKOV-ZIVANOV, N., TURNER, M. S., KANE, L. P., MOREL, P. A., AND FAEDER, J. R. The duration of t cell stimulation is a critical determinant of cell fate and plasticity. *Science Signaling* 6, 300 (Nov. 2013), ra97–ra97.
- [10] MISKOV-ZIVANOV, N., ZULIANI, P., CLARKE, E. M., AND FAEDER, J. R. Studies of biological networks with statistical model checking: application to immune system cells. In *Proceedings of the International Conference on Bioinformatics, Computational Biology and Biomedical Informatics*. ACM-DL, 2013, pp. 728–729.
- [11] MISKOV-ZIVANOV, N., ZULIANI, P., CLARKE, E. M., AND FAEDER, J. R. High-level modeling and verification of cellular signaling. In *18th IEEE International High Level Design Validation and Test Workshop (HLDVT)*. Institute of Electrical and Electronics Engineers, 2016, pp. 162–169.
- [12] SAYED, K., BOCAN, K. N., AND MISKOV-ZIVANOV, N. Automated extension of cell signaling models with genetic algorithm. In *2018 40th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)* (July 2018), IEEE.
- [13] SAYED, K., KUO, Y.-H., KULKARNI, A., AND MISKOV-ZIVANOV, N. Dish simulator: Capturing dynamics of cellular signaling with heterogeneous knowledge. In *2017 Winter Simulation Conference (WSC)* (2017), pp. 896–907.
- [14] SAYED, K., TELMER, C. A., BUTCHY, A. A., AND MISKOV-ZIVANOV, N. Recipes for translating big data machine reading to executable cellular signaling models. In *Lecture Notes in Computer Science*. Springer International Publishing, Dec. 2017, pp. 1–15.
- [15] VALENZUELA-ESCÁRCEGA, M. A., ÖZGÜN BABUR, HAHN-POWELL, G., BELL, D., HICKS, T., NORIEGA-ATALA, E., WANG, X., SURDEANU, M., DEMIR, E., AND MORRISON, C. T. Large-scale automated machine reading discovers new cancer-driving mechanisms. *Database* 2018 (Jan. 2018).
- [16] VAN DONGEN, S. *Graph clustering by flow simulation*. *Graph Stimulation by Flow Clustering*. PhD thesis, PhD thesis, University of Utrecht. <https://doi.org/10.1016/j.cosrev.2007....>, 2000.
- [17] WANG, Q., MISKOV-ZIVANOV, N., LIU, B., FAEDER, J. R., LOTZE, M., AND CLARKE, E. M. Formal modeling and analysis of pancreatic cancer microenvironment. In *Computational Methods in Systems Biology*. Springer International Publishing, 2016, pp. 289–305.
- [18] WHITLEY, D. A genetic algorithm tutorial. *Statistics and Computing* 4, 2 (June 1994).

Decodon Calculator 2: Codon-Optimized Degenerate Codon Set Design Tools

Akira Takada

The College of New Jersey
Ewing, New Jersey, USA

Tomer Aberbach

The College of New Jersey
Ewing, New Jersey, USA

Nicholas Carpino

The College of New Jersey
Ewing, New Jersey, USA

Georgios Papamichail

National Centre for Public
Administration
Athens, Greece

Dimitris Papamichail*

papamicd@tcnj.edu
The College of New Jersey
Ewing, New Jersey, USA

1 INTRODUCTION

Mutant libraries representing protein variants are often used to optimize protein function. One design strategy aims to place ambiguous codons at certain locations of a target protein, creating sequence degeneracy in the synthesized DNA library. Assays of such libraries aid the development of variants with improved properties [5, 6].

Combinatorial assembly of oligos with degeneracies provides a reasonable approach in the development of improved variants [1, 3, 8, 9]. Library-design strategies seek to experimentally evaluate a diverse but focused region of sequence space in order to improve the likelihood of finding a beneficial variant.

Such an approach is based on the premise that prior knowledge can inform generalized predictions of protein properties, but may not be sufficient to specify individual, optimal variants. Libraries are particularly appropriate when the prior knowledge does not admit detailed, robust modeling of the desired properties, but when experimental techniques are available to rapidly assay a pool of variants.

The design of mutant protein libraries typically involves selecting sites for mutation where degenerate codons (those containing mixtures of nucleotides) are introduced to enable variation. The protein variant library is then produced by synthesizing degenerate oligonucleotides and using annealing based recombination. Custom oligonucleotide overlaps enable the targeted introduction of crossovers at only desired positions, in turn enabling the desired level and type of diversity in the combinatorial library [1, 3, 8, 9].

Traditional mutant protein library design methods involve the incorporation of a single degenerate codon (thereafter referred to as *decodon*) at each position where amino acid substitutions are considered. Decodons contain ambiguous (*degenerate*) bases, as shown in Table 1. Degenerate bases are one letter codes are used to represent (i.e. code) sets of DNA bases.

Table 1: Degenerate Bases and their codings

Degenerate Base	Actual Bases Coded
N	A or C or G or T
B	C or G or T
D	A or G or T
H	A or C or T
V	A or C or G
K	G or T
M	A or C
R	A or G
S	C or G
W	A or T
Y	C or T

CodonGenie [7] is an online tool that was created to aid the effort of designing single decodons that code for any given amino acid set. The tool ranks candidate decodons by specificity, attempting to minimize coding of undesirable amino acids and/or STOP codons. Even so, when using a single decodon to code for a set of amino acids, it is often unavoidable to code for additional unwanted amino acids and/or STOP codons.

2 A CODON-OPTIMIZED DECODON CALCULATOR TOOL

In our previous work [4] we explored coding sets of amino acids by using multiple decodons. Annealing-based recombination of degenerate oligos containing these decodons can generate libraries that focus exclusively on the productive portion of the design space by eliminating unwanted variants, therefore improving the overall quality of the library for screening purposes. This first version of our Decodon Calculator, given as input any set of amino acids, produces the minimum number of decodons necessary to code for exactly that set, i.e. without coding for extraneous amino

*Corresponding author

acids or STOP codons. At the same time it outputs an example of a decodon set of minimum cardinality for each amino acid subset. These example sets are produced based on the order that decodon sets are generated and processed by our dynamic programming algorithm and, for all purposes, can be considered random. To enhance the utility of our tool, in the second generation Decodon Calculator we introduce a decodon set scoring scheme, based on individual organism codon preferences.

To score decodon sets in a given organism we utilize the codon usage frequency tables in the Codon Usage Database [2]. For each codon, we calculate the ratio of the frequency of that codon over the frequency of the most frequent synonymous codon in the organism. Then, to score a set of decodons, we calculate the numeric average of all individual codons comprising the set. Therefore scores for decodon sets range from 0 - 1, with a score of 1 indicating that each codon in the set is the most frequently occurring in the organism of interest. This scoring method is similar to the one utilized by the CodonGenie tool, and more details can be found in [7].

The new Decodon Calculator, given any input amino acid set, now outputs the minimum cardinality decodon set with the largest codon utilization (referred in the tool as *organism compatibility*) score coding for these amino acids. Currently our tool supports codon preferences for three organisms, human, mouse, and E.coli. Additional model organisms are expected to be added in the future.

3 SUBSETS AND SUPERSETS

The Decodon Calculator provides minimum cardinality decodon sets coding for any given amino acid set, eliminating unwanted mutations. When designing oligos though with multiple mutation sites, the degeneracy can increase exponentially as a function of the number of sites. This in turn, when the desired number of mutations per oligo is large, can lead to substantially increases in the total number of distinct oligos that need to be ordered. To balance DNA synthesis costs and library specificity, we added an additional feature to the Decodon Calculator, which provides maximal subsets and minimal supersets of the input amino acid sets which are encoded by fewer decodons than the original input set necessitates.

In addition to the minimal cardinality k decodon set, our tool now also generates the maximal subset and minimal superset of the input amino acid set encoded by $k-1$ decodons, then by $k-2$, etc. The process stops once we reach subsets/supersets encoded by a single decodon, or when the cardinality difference exceeds 3, meaning the subset/superset has more than 3 fewer/additional amino acids than the original input set.

The feature works as follows: once a set of amino acids is selected and the *Submit* button is pressed, a table appears

Optimal Degenerate Codon Design for Amino Acid Sets for Specific Organisms

By Akira Takada & Dimitris Papamichail @ [The College of New Jersey \(TCNJ\)](http://The College of New Jersey (TCNJ))

Select a set of amino acids. Select an organism. Click on 'Submit' to calculate the most optimal set of degenerate codons needed to code for the amino acids for the organism.

Non-polar
 A F Glycine (Gly) I L M Proline (Pro) V Tryptophan (Try)

Polar
 C Asparagine (Asn) Q Serine (Ser) T Y

Acidic
 Aspartic Acid (Asp) E

Basic
 H Lysine (Lys) R

E.Coli

Minimum Number of Degenerate Codons: 4
 Degenerate Codon(s) Example: AAA CCG RRT TGG
 Compatibility Score: 0.995 / 1

# of Degenerate Codon(s)	Optimal Subset	Subset	Score	Optimal Superset	Subset	Score
3	P,N,S,D,G,W	CCG RRT TGG	0.994 / 1	K,N,D,G,P,R,S,W	AAW GRT YSG	0.922 / 1
2	N,S,D,G,W	RRT TGG	0.991 / 1	K,N,R,S,E,D,G,P,W	RRW YSG	0.818 / 1
1	N,S,D,G	RRT	0.981 / 1	NULL	NULL	NULL

Figure 1: Example result of Decodon Calculator showing minimal decodon set with maximal codon score and corresponding optimal subsets and supersets.

on the bottom of the screen, as shown in in Figure 1. In this particular example, we can observe that the residues G, P, N, S, D, and K can be coded by the four decodons AAA, CCG, RRT, and TGG. Using three decodons, it is possible to code for six of the seven input amino acids, P, N, S, D, G, and W. Or for the superset {K, N, D, G, P, R, S, W}, which includes one additional amino acid. With one decodon it is possible to code for four of the seven input amino acids, but a superset with at most 3 additional residues cannot be encoded by a single decodon.

The Decodon Calculator featuring organism optimal decodons can be accessed at <http://algo.tcnj.edu/decodoncalc2/>. The Decodon Calculator with the additional optimal subset and superset functionality can be accessed at <http://algo.tcnj.edu/decodoncalcset/>. Both tools have been implemented in Javascript on the client side and PHP on

the server side, and utilize a database to store and retrieve optimal decodon sets for all amino acid set.

4 ACKNOWLEDGEMENTS

The authors acknowledge use of the ELSA high performance computing cluster at The College of New Jersey for hosting the web service reported in this paper. This cluster is funded by the NSF grant OAC-1828163.

REFERENCES

- [1] MEYER, M. M., SILBERG, J. J., VOIGT, C. A., ENDELMAN, J. B., MAYO, S. L., WANG, Z.-G., AND ARNOLD, F. H. Library analysis of SCHEMA-guided protein recombination. *Protein Science* (2003).
- [2] NAKAMURA, Y., GOJOBORI, T., AND IKEMURA, T. Codon usage tabulated from international DNA sequence databases: status for the year 2000. *Nucleic Acids Research* 28, 1 (01 2000), 292–292.
- [3] PANTAZES, R. J., SARAF, M. C., AND MARANAS, C. D. Optimal protein library design using recombination or point mutations based on sequence-based scoring functions. *Protein Engineering, Design and Selection* (2007).
- [4] PAPAMICHAIL, D., CARPINO, N., ABERBACH, T., AND PAPAMICHAIL, G. Decodon Calculator: Degenerate Codon Set Design for Protein Variant Libraries. *12th International Workshop on Bio-Design Automation* (2020).
- [5] PARKER, A. S., ZHENG, W., GRISWOLD, K. E., AND BAILEY-KELLOGG, C. Optimization algorithms for functional deimmunization of therapeutic proteins. *BMC Bioinformatics* (2010).
- [6] REETZ, M. T., AND CARBALLEIRA, J. D. Iterative saturation mutagenesis (ISM) for rapid directed evolution of functional enzymes. *Nature Protocols* (2007).
- [7] SWAINSTON, N., CURRIN, A., GREEN, L., BREITLING, R., DAY, P. J., AND KELL, D. B. CodonGenie: Optimised ambiguous codon design tools. *PeerJ Computer Science* (2017).
- [8] TREYNOR, T. P., VIZCARRA, C. L., NEDELCO, D., AND MAYO, S. L. Computationally designed libraries of fluorescent proteins evaluated by preservation and diversity of function. *Proceedings of the National Academy of Sciences of the United States of America* (2007).
- [9] VOIGT, C. A., MARTINEZ, C., WANG, Z. G., MAYO, S. L., AND ARNOLD, F. H. Protein building blocks preserved by recombination. *Nature Structural Biology* (2002).

Data Representation in the DARPA SD2 Program

Nicholas Roehner¹, Jacob Beal¹, Bryan Bartley¹, Richard Markeloff¹, Tom Mitchell¹, Tramy Nguyen¹, Daniel Sumorok¹, Nicholas Walczak¹, Chris Myers², Zach Zundel³, James Scholz³, Benjamin Hatch³, Mark Weston⁴, John Colonna-Romano⁵

¹Raytheon BBN Technologies, ²University of Colorado Boulder, ³University of Utah, ⁴Netrias, ⁵Aptima
nicholas.roehner@raytheon.com, jakebeal@ieee.org, chmy5075@colorado.edu, weston@netrias.com

1 MOTIVATION

Modern scientific enterprises are often highly complex and multidisciplinary, particularly in areas like synthetic biology where the subject at hand is itself inherently complex and multidisciplinary. Collaboration across many organizations is necessary to efficiently tackle such problems [6, 15], but remains difficult. The challenge is further amplified by automation that increases the pace at which new information can be produced, and particularly so for matters of fundamental research, where concepts and definitions are inherently fluid and may rapidly change as an investigation evolves [7].

The DARPA program Synergistic Discovery and Design (SD2) aimed to address these challenges by organizing the development of data-driven methods to accelerate discovery and improve design robustness, with one of the key domains under study being synthetic biology. The program was specifically organized such that teams provided complementary types of expertise and resources, and without any team being in a dominant organizational position, such that subject-matter investigations would necessarily require peer-level collaboration across multiple team boundaries. With more than 100 researchers across more than 20 organizations, several of which ran experimental facilities with high-throughput automation, participants were forced to confront challenges around effective data sharing.

The default architecture for scientific collaboration is essentially one of anarchy, with ad-hoc bilateral relations between pairs of collaborators or experimental phases (Figure 1(a)). This was by necessity the case during early phases of the SD2 program as well, in which incorporating new tools into pipelines was ad-hoc and time-consuming, and data was generally disconnected from genetic designs and experimental plans. The other typical approach for collaboration is one of “command and control”, in which a dominant organization determines the data sharing content and format for all participants (Figure 1(b)). This can be efficient, but tends to be limited in flexibility and extensibility, rendering it unsuitable for research collaboration, as indeed was found when we attempted this approach during the first year of the SD2 program. We addressed these problems with the application of distributed standards to create a “flexible rendezvous” model of collaboration (Figure 1(c)), enabling information flow to track evolving collaborative relationships, improving

the sharing and utility of information across the community and supporting accelerated rates of experimentation.

2 APPROACH

The driving design philosophy behind our approach to user interaction in SD2 was to adapt representational tooling as closely as possible to existing tools and familiar interfaces, such as spreadsheets and word processing documents. Taking this approach allowed us to use and improve formal machine-readable representations for system integration while minimizing the amount that participating researchers needed to learn about the formal representations. The central set of standards thus formed a point of rendezvous between the various stakeholders interacting in different experimental roles, while still allowing each of these participants to continue working in their native idiom.

Specifically, the data sharing working group collaborated in the creation of several key advancements in standards and tooling that combine to create a comprehensive ecosystem of lightweight curation tools. These are:

- Advances in biological data representation in the form of enhancements to SBOL2 for comprehensive representation of the design-build-test-learn cycle [5, 8] and ultimately the development of SBOL3 [11], which in turn enabled us to accelerate the rate of data standard development and integration.
- A plugin interface to SynBioHub [12], which enabled rapid development of new functionality for data visualization, submission, and exchange [10].
- The SBOL Project Dictionary tool [4], which provides a Google Sheets interface for collective “just-in-time” harmonization of terminology across organizations, thus enabling metadata translation and data fusion.
- The Experimental Intent Parser tool [13], an extension of Google Docs that enabled biologists to design and launch automated experiments with an easy-to-use interface.
- The Open Protocol Interface Language (OPIL) [1], which enabled the laboratories executing experiments to share information about their protocols with experiment planning tools, thus better informing the investigators proposing experiments to execute.

- The SYNBICT tool [14], which enabled automated generation of improved annotations and extraction of functional models of biological designs from their sequences.
- The REDOER tool [2], which attempts to infer experimental design from collections of samples, enabling quality-control on automated experimentation.
- The Excel2SBOL conversion tool [9], which enabled an efficient workflow for producing build requests for genetic designs.

Collectively deployed in the architecture shown in Figure 2, these tools enabled a shift in the organization of experimental and informational workflows toward faster and more flexible execution, most notably in SD2 working groups that were working on challenge problems focused on the performance of genetic circuits in yeast, moving existing designs into novel chassis, and cell-free riboswitch design.

Following this shift, program knowledge sharing expanded greatly. One key measure of knowledge sharing is the number of terms stored in the SBOL Project Dictionary, as each such term indicates a strain, reagent, genetic construct, parameter, or other similar item that is being communicated between collaborating organizations. We find that the number of terms stored in the SBOL Project Dictionary, expanded in close correlation with the increase in experiment tempo: Figure 3 shows the correlation between knowledge sharing and data production, with both moving much more quickly in the second half of SD2 after these tools began to be released.

Moreover, measurement of key knowledge collections shows that progress on tools correlates with increases in knowledge sharing and productivity. Figure 4 shows that knowledge expansions in SynBioHub are correlated with the dates of tool releases. In particular, key tool releases occurred around July 2018 (Project Dictionary), April 2019 (SYNBICT, REDOER, and Intent Parser), October 2019 (Experiment launches via Intent Parser), January 2020 (Excel2SBOL), April 2020 (SYNBICT libraries), and October 2020 (OPIL), and these are correlated with expansions in the number of SBOL Module relations, which are used in representing knowledge about circuits, and the number of ModuleDefinition relations, which are used in representing strains and reagents.

Finally, beyond these specific knowledge classes, the overall volume of knowledge systematized by the program is quite large as well. By the end of the period reported, the SD2 SynBioHub instance had a knowledge store of 22,872,306 triples, including 3,549,237 components, 15,884 modules, and 524,377 collections.

3 DISCUSSION

The transition partners that we have engaged with these tools see potential value in using the Synthetic Biology Open

Language (SBOL) and associated tools with well-defined APIs as a standard that is relevant to many groups and has greater potential for sharing, greater levels of support, and more longevity. Partners can benefit from work by the supporting community and worry less about their data-sharing infrastructure losing “product support.” Transition partners also see value in standards and tools being free and open, which allows them to know what their data-sharing methods are doing, modify them if needed, interact with them easily via their own code, and not be limited by high commercial costs or number of licensed “seats.” With solutions like those that we developed in the SD2 program, they will be better able to keep track of information over time as new people come and go from labs, so that they can continue to build new knowledge on top of existing knowledge. They will also be better able to share design information and associated data with other research groups in a more consistent way, and be better able to take advantage of other groups’ designs and data, making their own engineering processes much faster.

Building on the success of data representation in the SD2 program, we recommend that future programs should support the development of standards and corresponding software infrastructure and also should support curators and the development of data repositories and curation tools. Similarly, just as many government funding agencies now require open access publications, government funding agencies should require funded science activities to use standard-enabled workflows and software to enhance data management and sharing, and should ensure that funding is specifically allocated for such activities. Looking farther ahead, we also see an opportunity for increased machine-readability of shared information to become a foundation for higher-level autonomy in scientific investigation [3], as well as for enhancing reproducibility through machine validation of scientific experiments and automation-assisted publication of experiments.

4 ACKNOWLEDGEMENTS

This work was supported by Air Force Research Laboratory (AFRL) and DARPA contracts FA8750-17-C-0184, FA8750-17-C-0229, FA8750-17-C-0231, and FA8750-17-C-0294. This document does not contain technology or technical data controlled under either U.S. International Traffic in Arms Regulation or U.S. Export Administration Regulations. Views, opinions, and/or findings expressed are those of the author(s) and should not be interpreted as representing the official views or policies of the Department of Defense or the U.S. Government.

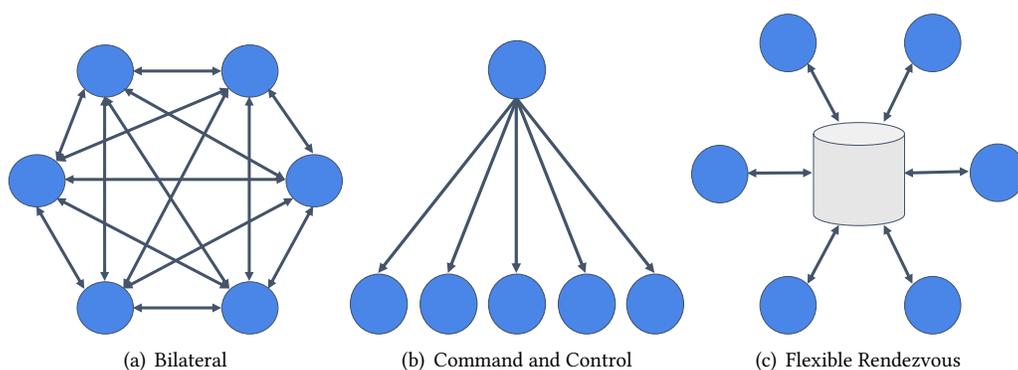


Figure 1: Architectures for data sharing: bilateral relations (a), command and control (b), and flexible rendezvous (c).

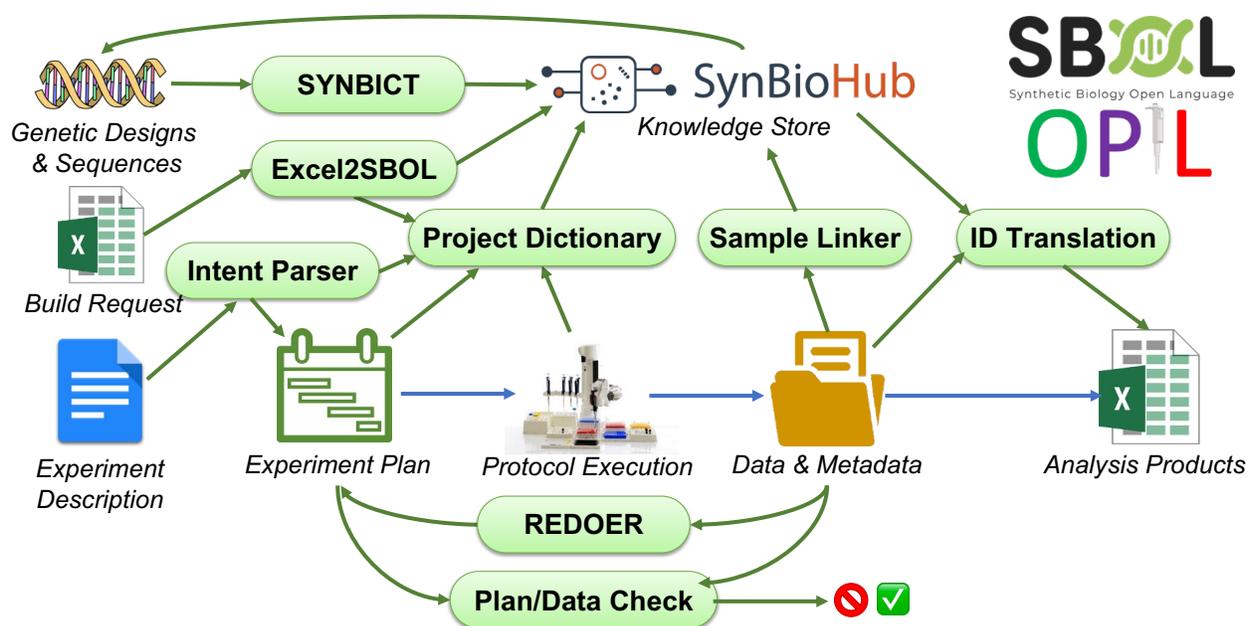


Figure 2: High-level diagram showing how data representation tools were deployed in the DARPA SD2 program with respect to inputs of designs, build requests, and experiments requests, and outputs of data, metadata, and analysis products. This diagram focuses specifically on the key representations (SBOL and OPIL) and representation-centric tooling, and not other aspects of supporting automation used in SD2.

REFERENCES

- [1] BARTLEY, B., BEAL, J., BRYCE, D., GOLDMAN, R. P., KELLER, B., LADWIG, J., LEE, P., MARKELOFF, R., NGUYEN, T., NOWAK, J., AND WESTON, M. Open protocol interface language. <https://github.com/SD2E/OPIL-specification>, 2021.
- [2] BARTLEY, B., BEAL, J., AND WESTON, M. Reverse engineering design of experiments for review (redoer). In *AI4SynBio at AAAI SSS* (March 2021).
- [3] BEAL, J., AND ROGERS, M. Levels of autonomy in synthetic biology engineering. *Molecular Systems Biology* 16, 12 (2020), e10019.
- [4] BEAL, J., SUMOROK, D., BARTLEY, B., AND NGUYEN, T. Collaborative terminology: Sbol project dictionary. In *12th International Workshop on Bio-Design Automation (IWBDA)* (August 2020).
- [5] COX, R. S., MADSEN, C., MCLAUGHLIN, J. A., NGUYEN, T., ROEHNER, N., BARTLEY, B., BEAL, J., BISSELL, M., CHOI, K., CLANCY, K., ET AL. Synthetic biology open language (sbol) version 2.2. *Journal of integrative bioinformatics* 15, 1 (2018).
- [6] LAKHANI, K. R., AND PANETTA, J. A. The principles of distributed innovation. *Innovations: technology, governance, globalization* 2, 3 (2007), 97–112.
- [7] LATOUR, B., AND WOOLGAR, S. *Laboratory life*. Princeton University Press, 2013.
- [8] MADSEN, C., MORENO, A. G., UMESH, P., PALCHICK, Z., ROEHNER, N., ATALLAH, C., BARTLEY, B., CHOI, K., COX, R. S., GOROCHOWSKI, T., ET AL. Synthetic biology open language (sbol) version 2.3. *Journal of*

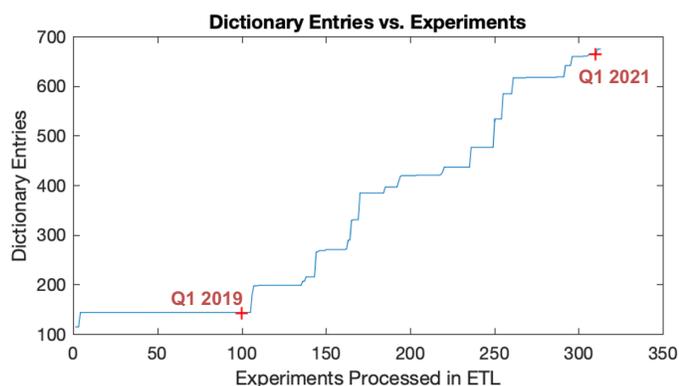


Figure 3: Increased knowledge sharing is correlated with overall rates of experimentation in SD2. Accumulation of shared knowledge, as measured by increased numbers of entries in the SBOL Project Dictionary, increased much more rapidly in the second half of the SD2 program, as did the rate at which experiments were run.

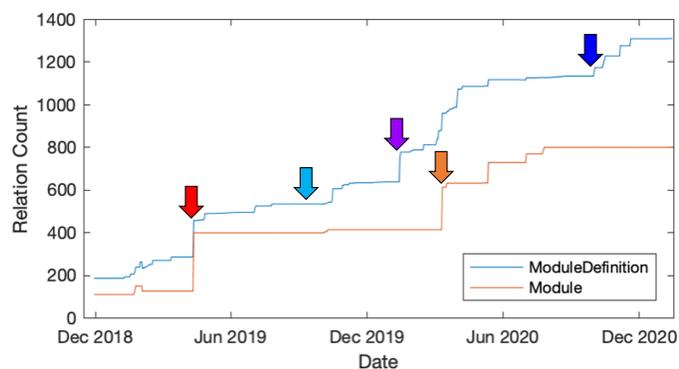


Figure 4: Introduction of specific tools in SD2 is correlated with expansions in key forms of shared knowledge. Arrows mark the approximate time of introduction of key advances in data representation tools: SYNBICT and Experimental Intent Parser (red), experiment launches via Experimental Intent Parser (light blue), Excel-to-SBOL (purple), SYNBICT libraries (orange), and OPIL (dark blue).

a standards-enabled design repository for synthetic biology. *ACS synthetic biology* 7, 2 (2018), 682–688.

- [13] NGUYEN, T., WALCZAK, N., BEAL, J., SUMOROK, D., AND WESTON, M. Intent parser: a tool for codifying experiment design. In *12th International Workshop on Bio-Design Automation (IWBDA)* (August 2020).
- [14] ROEHNER, N., MANTE, J., MYERS, C. J., AND BEAL, J. Synthetic biology curation tools (synbict). *under review*.
- [15] VON HIPPEL, E. “sticky information” and the locus of problem solving: implications for innovation. *Management science* 40, 4 (1994), 429–439.

integrative bioinformatics 16, 2 (2019).

- [9] MANTE, J., POTZSCH, I., ABAM, J., BEAL, J., AND MYERS, C. J. Excel-sbol converter: Creating sbol from excel templates and vice versa. In *Submitted to 13th International Workshop on Bio-Design Automation (IWBDA)* (September 2021).
- [10] MANTE, J., ZUNDEL, Z., AND MYERS, C. Extending synbiohub’s functionality with plugins. *ACS synthetic biology* 9, 5 (2020), 1216–1220.
- [11] McLAUGHLIN, J. A., BEAL, J., MISIRLI, G., GRÜNBERG, R., BARTLEY, B. A., SCOTT-BROWN, J., VAIDYANATHAN, P., FONTANARROSA, P., OBERORTNER, E., WIPAT, A., ET AL. The synthetic biology open language (sbol) version 3: simplified data exchange for bioengineering. *Frontiers in Bioengineering and Biotechnology* 8 (2020), 1009.
- [12] McLAUGHLIN, J. A., MYERS, C. J., ZUNDEL, Z., MISIRLI, G., ZHANG, M., OFITERU, I. D., GONI-MORENO, A., AND WIPAT, A. Synbiohub:

Network visualisation of synthetic biology designs

Matthew Crowther
m.crowther1@ncl.ac.uk
Newcastle University
Newcastle Upon Tyne, United Kingdom

Anil Wipat
anil.wipat@ncl.ac.uk
Newcastle University
Newcastle Upon Tyne, United Kingdom

Ángel Goñi-Moreno
angel.goni@upm.es
Universidad Politécnica de Madrid
Madrid, Spain

1 ABSTRACT

Visualising the complex information captured by synthetic biology designs is still a major challenge. The popular glyph approach where each genetic part is displayed on a linear sequence allows researchers to generate diagrams and visualise abstract designs [2], but only represents a single, static representation that results in visualisation that is not specific to the requirements of a user resulting in a one-size-fits-all visualisation. We developed a network visualisation technique that automatically turns all design information into a graph, displaying otherwise hidden data. The structure of the resulting graphs can be dynamically adjusted according to specific visualisation requirements, such as highlighting proteins, interactions or hierarchy. Since biological systems have an inherent affinity with network visualization [6], we advocate for adopting this approach to standardise and automate the representation of complex information.

2 RESULTS

Firstly, a NOR gate design (adapted from [8]) is used to showcase some fundamental visualization processes and the methods. Secondly, more complex regulatory circuits are used to illustrate the potential of the network visualisation approach to effectively display novel features.

Data. Before any visualisation can be realised, the underlying data representation must be considered. Without a rich data representation, most meaningful visualisation is not achievable simply due to the data not being encoded. Therefore, we will use the Synthetic Biology Open Language (SBOL) [5] as the data capture format. SBOL is a more formal and synthetic biology-centric approach to the design specification.

View. Visualizing unmodified data will produce an incomprehensible visualisation as the domain is too broad so the significance of connections is lost. As seen within Figure 1A, despite a small design due to the verbose nature of the underlying data very little can be inferred. Therefore, a view is defined as an aggregation of data to produce a graph that is focused on a specific aspect of the design. With a more concentrated domain focusing on a single aspect, visual complexity is reduced. For this introduction, a basic view that

aggregates data into the overall design and constituent biological parts and entities is used as seen within Figure 1B.

Layout. A view of any meaningful size will produce an incoherent visualization when the position of nodes do not consider the data being represented. Layout pertains to the coordinate location of the nodes within the plot. Trivially, layouts can ensure the rendered nodes and edges do not overlap but more significant implementation can have a layout mirroring the intent behind the data being visualized. In the current working example (Figure 1A, despite a better visual, the network is incoherent as no positional data is encoded). However, as seen within Figure 1C, despite a relatively basic layout, the visual output is more clear.

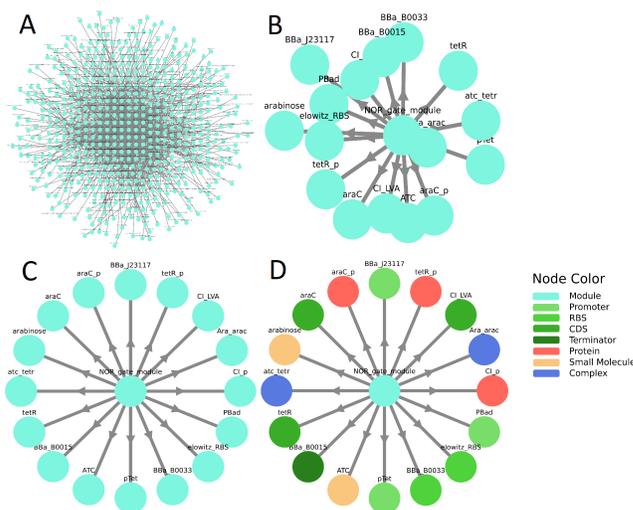


Figure 1: A) All encoded data is rendered. B) Simple view of constituent biological parts, proteins and non-genetic entities described within a design. C) A basic concentric layout - all constituent entities are positioned around the central node denoting the overall design. D) Addition of colour mapping with biological entities types and roles within the design.

Label reduction. Labels are added directly to the graph, connected to edges or nodes. However, screen space is finite and can become saturated. This issue is compounded when: nodes are closely positioned, the rendered text is long, and/or

the graph is highly connected. Therefore, label reduction is the process of replacing labels with visual features to increase concision but still encode the information. With the working example, while the core focus on the view is comprehensible, information that may be desired is not present. As seen within Figure 1D, a user may want to visualise the role of each biological entity.

Visualising complex information via *presets*

The overview previously discussed is only one instance of producing a comprehensible visualisation. Here we use the term “preset” to denote a view combined with a collection of visual techniques that are complementary to said view such that the visual output focuses attention upon a specific and desired feature of a design. Below two presets are discussed, including intent and how visual modifications have an affinity with the view. However, this is not an exhaustive list and providing the information is encoded within the design data, any feature of a design can be visualised using a network/graph focused approach.

Hierarchy. The hierarchy of a design focuses on visualizing how the different perceived levels of biological entities are structured. This provides insight into each abstraction level and how the components of each level map to their neighbours. Furthermore, a hierarchical view can visualize a design of arbitrary depth which is beneficial since the levels of abstraction within a design increase as modules become larger (parts, devices, circuits, systems, consortia). Figure 2 displays a hierarchical view that allows not only a visualization of individual parts and constructs but also the makeup of larger modules.

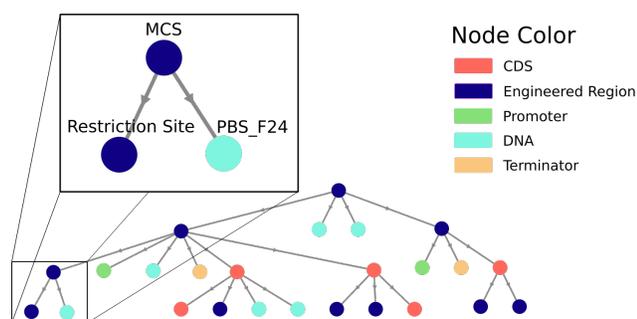


Figure 2: Visualising the *digitalizer* synthetic circuit from [3] with a hierarchical representation of genetic modules in 4 layers.

Interaction. In contrast with sequence-level visualisations, where intent and function are not explicitly described and non-genetic entities are often poorly represented, interaction networks provide an explicitly functional perspective. By

using this view, non-genetic entities (e.g., proteins) are easily represented. Furthermore, the description of biochemical networks fit well into a graph-based approach since these are conceptualized as a set of interacting entities. We visualised a Boolean genetic circuit (Figure 3) by using only its interactions and non-genetic elements. Inputs, outputs and information flow are easily comprehended even at a glance. However, visualising the same relatively complex design using sequence-level information by deriving functional details would be more challenging.

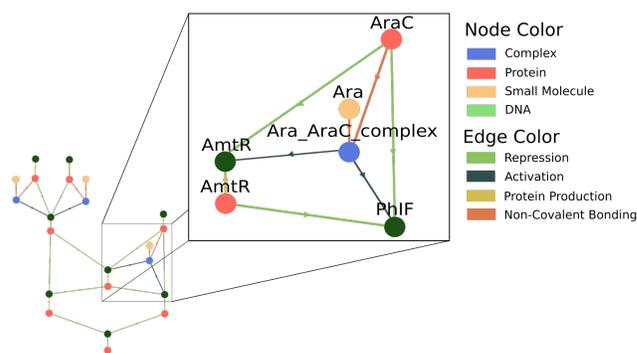


Figure 3: Visualising interactions (edges) and non-genetic elements (nodes) of the *0xC7* Boolean genetic circuit from [7].

Scaling abstraction. Very often, despite the focus on a specific design feature (e.g., non-genetic elements), issues of comprehension still arise due to the level of design details and annotations. The ability to visualize a higher level of abstraction allows a more granular, and more easily comprehensible output [4]. Figure 4 displays the same design as in Figure 3 but at a higher level of abstraction, which may be more adequate for a rapid design check.

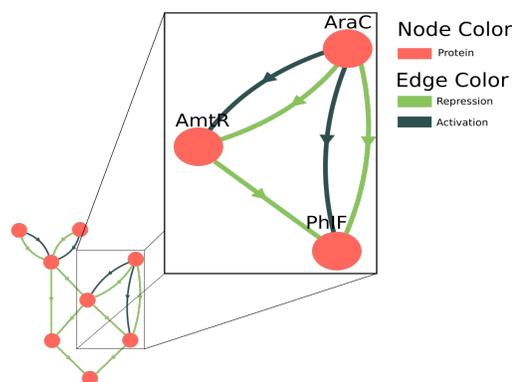


Figure 4: Visualising interactions (edges) and non-genetic elements (nodes) of the *0xC7* Boolean genetic circuit from [7] with higher abstraction.

While increasing the abstraction level can produce a visual output that is more comprehensible in terms of function, the reduced granularity can lead to more ambiguous visualisations concerning mechanistic details. Therefore, lowering the level of abstraction (thus visualising more details) may be beneficial in some cases, for instance, when building a mathematical model of a genetic design. Figure 5 is a more detailed view compared to Figure 3 and despite a considerable complexity increase, interactions are broken down into a number reactions thus providing more detailed information.

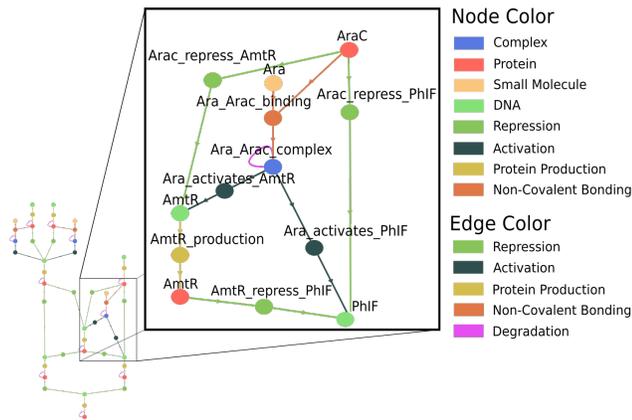


Figure 5: Visualising interactions (edges) and non-genetic elements (nodes) of the 0xC7 Boolean genetic circuit from [7] with lower abstraction.

3 DISCUSSION & FUTURE WORK

We present a visualisation method that offers an alternative approach to a conventional genetic-parts-based glyph. Designs can be automatically produced at differing levels of abstraction, as defined by the user. This approach promises to help understand complex designs more easily, and to scale better for large designs, such as chromosomes.

Future efforts will focus on four aspects. Firstly, the fact that networks are generated automatically and adjusted dynamically paves the way to develop powerful user interaction tools. Secondly, we will exploit the full potential of graphs for mathematical analysis through a wealth of graph theory methods. Indeed, networks are not only useful for visualisation purposes but mathematical structures for studying data. Thirdly, networks allow representing any type of data, not just gene design information. Therefore, specific visualisation networks of every stage throughout a standardised DBTL lifecycle [1] will be coupled into layered graphs that will include from automation to characterisation to modelling information. Finally, most current visualisation techniques will not scale to large designs. Therefore, exploring

how network visualisation that has precedence with large data visualisation can be applied to designs of extreme size.

Visualization is complementary to the development of data standards. Here, we use designs encoded with the SBOL since this captures richer information than GenBank or FASTA formats.

REFERENCES

- [1] BEAL, J., GOÑI-MORENO, A., MYERS, C., HECHT, A., DE VICENTE, M. D. C., PARCO, M., SCHMIDT, M., TIMMIS, K., BALDWIN, G., FRIEDRICH, S., ET AL. The long journey towards standards for engineering biosystems: Are the molecular biology and the biotech communities ready to standardise? *EMBO reports* 21, 5 (2020), e50521.
- [2] BEAL, J., NGUYEN, T., GOROCHOWSKI, T. E., GOÑI-MORENO, A., SCOTT-BROWN, J., McLAUGHLIN, J. A., MADSEN, C., ALERITSCH, B., BARTLEY, B., BHAKTA, S., ET AL. Communicating structure and function in synthetic biology diagrams. *ACS synthetic biology* 8, 8 (2019), 1818–1825.
- [3] CALLES, B., GOÑI-MORENO, Á., AND DE LORENZO, V. Digitalizing heterologous gene expression in gram-negative bacteria with a portable on/off module. *bioRxiv* (2019).
- [4] HEINEMANN, M., AND PANKE, S. Synthetic biology—putting engineering into biology. *Bioinformatics* 22, 22 (09 2006), 2790–2799.
- [5] MADSEN, C., MORENO, A. G., UMESH, P., PALCHICK, Z., ROEHNER, N., ATALLAH, C., BARTLEY, B., CHOI, K., COX, R. S., GOROCHOWSKI, T., ET AL. Synthetic biology open language (sbol) version 2.3. *Journal of integrative bioinformatics* 16, 2 (2019).
- [6] MIELE, V., MATIAS, C., ROBIN, S., AND DRAY, S. Nine quick tips for analyzing network data. *CoRR abs/1904.05334* (2019).
- [7] NIELSEN, A. A. K., DER, B. S., SHIN, J., VAIDYANATHAN, P., PARALANOV, V., STRYCHALSKI, E. A., ROSS, D., DENSMORE, D., AND VOIGT, C. A. Genetic circuit design automation. *Science* 352, 6281 (2016).
- [8] TAMSIR, A., TABOR, J. J., AND VOIGT, C. A. Robust multicellular computing using genetically encoded nor gates and chemical ‘wires’. *Nature* 469, 7329 (Jan 2011), 212–215.

A database for prokaryotic ligand-inducible transcription factors

Simon d'Oelsnitz

Danny Diaz*

doelsnitz@utexas.edu

The University of Texas at Austin

Austin, Texas

1 INTRODUCTION

Innovation in analytical chemistry is critical for the growing \$14.1 B high-throughput screening industry. While traditional chromatography and mass spectrometry methods are commonplace for chemical measurement, they often suffer numerous limitations [1]. (1) The separation of chemical species via chromatography is typically too slow for high-throughput screening applications, requiring 5-30 minutes per sample; while mass spectrometry alone offers rapid analyses, on the order of 1 second per sample, compounds with identical exact masses cannot be distinguished. (2) Target compounds lacking chromophores require derivatization, which increases analysis time and operational cost. (3) Analysis in complex backgrounds may produce a matrix effect, which reduces measurement performance. (4) Chromatographic and mass spectrometry instruments are often expensive and require specialized training to operate.

Nature has evolved elegant tools to bypass these limitations - genetic biosensors. In particular we focus on one-component prokaryotic ligand-inducible transcription factors (herein called "biosensors"), which have a history of use in synthetic genetic systems to report on the abundance of diverse chemical structures in a wide range of host organisms. [2]. Alternative sensor families (such as RNA-based sensors, antibodies, eukaryotic sensors, and two-component prokaryotic systems) are less commonly used as chemical measurement tools by comparison. Briefly, biosensors function in these systems by binding to a specific DNA sequence, thereby preventing the expression of the green fluorescent protein (GFP) gene, in the absence of the inducer molecule. In the presence of the target molecule, the biosensor dissociates from its DNA sequence, enabling GFP expression (Figure 1). As a result, the concentration of a particular chemical species is transduced into an easily-detectable fluorescent signal.

Biosensors offer numerous benefits over traditional chromatographic and mass spectrometry methods. Protein-chemical interactions employed by biosensors enable high specificity, even among enantiomers, which obviates derivatization and bypasses matrix effects. Furthermore, analysis occurs at the single-cell level, enabling the facile measurement of hundreds of thousands of samples in parallel.

Despite these clear advantages genetic biosensors are rarely used for chemical analysis, partly due to the fact that information on genetic biosensors is scattered across thousands of academic journals. To facilitate the accessibility of information of biosensors, we have created The Genetic Biosensor Database, which documents key features of biosensors extracted from academic journals, including metadata, structural data, the sensor's DNA binding sequence and cognate ligand, and appropriate references. This resource aims to facilitate the development of genetic biosensor systems for next generation chemical analysis.

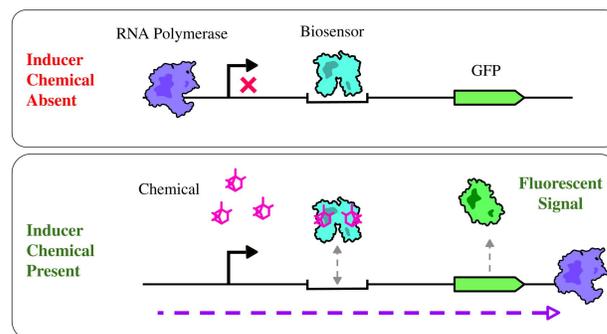


Figure 1: Mechanism of repressor-based prokaryotic ligand-inducible transcription factors.

2 THE GENETIC BIOSENSOR DATABASE

The objective of this database is to provide all information necessary to create a genetic reporter system (Figure 1). This includes (1) the sensor's cognate ligand(s), (2) the sensor's corresponding DNA sequence, and (3) the sensor's protein sequence. In addition, metadata, genome context, and structural data is provided (if available) to describe the sensor's role in its native host organism and possibly to infer unannotated functions. References are included for all data presented.

The database organizes each biosensor into one of eight transcription factor families: TetR, LysR, AraC, MarR, LacI, GntR, LuxR, and IclR (Figure 2). Once a family is selected, the

user may browse through a table of sensors belonging to that family. After a sensor is selected, data specific for that sensor appears below. The top section of information includes the sensor's alias, or NCBI accession number, a toggleable panel displaying the chemical structure of the sensor's cognate ligand(s), and a table containing corresponding metadata.

The middle section displays the sensor's cognate DNA binding sequence, or operator, and an illustration of the genome context of the sensor. Neighboring genes are color coded according to their annotated function; hovering over the gene displays its annotation and clicking on it leads to the corresponding NCBI entry page.

The final section below displays the sensor protein's sequence as well as a toggleable and interactive display of its structure(s), if available. Finally, the bottom of each page contains appropriate references with links to the corresponding article and labels describing what information the reference provides; this can include data on a ligand interaction, DNA interaction, structure, or an application using the sensor.

Data curation

All data is manually curated from peer-reviewed academic publications. To encourage the development of functional genetic reporter systems, a relatively strict criteria for including information on sensor ligand/DNA interactions is enforced. To be included in this database, data supporting a ligand-sensor interaction must be derived from an EMSA, SPR, or ITC *in vitro* experiment or a synthetic *in vivo* reporter system within a heterologous host (such as in Figure 1). *In vivo* data reporting on mRNA or protein abundance in the native host organism upon exposure to the target chemical is not sufficient for inclusion in the database, since this assay cannot distinguish whether the input chemical or a metabolic intermediate produced by native enzymes is the true inducer molecule for the sensor [3]. Additionally, structures of the sensor protein in complex with the inducer molecule is not sufficient to establish a ligand-sensor interaction, since crystallography conditions can produce biases and there have been documented cases of a sensor being crystallized in complex with non-inducer molecules [4].

Data supporting a ligand-DNA interaction must be derived from an EMSA or DNase footprinting *in vitro* experiment, or a synthetic *in vivo* reporter system within a heterologous host (such as in Figure 1), to be included in the database. These sequences typically include inverted repeats. In the case where a sensor binds to multiple different DNA sequences, the sequence with the highest binding affinity is chosen.

Future work will include (A) integrating full-text search, (B) automated curation of biosensor entries, and (C) implementing tools for predicting a sensor's operator sequence and cognate ligands.

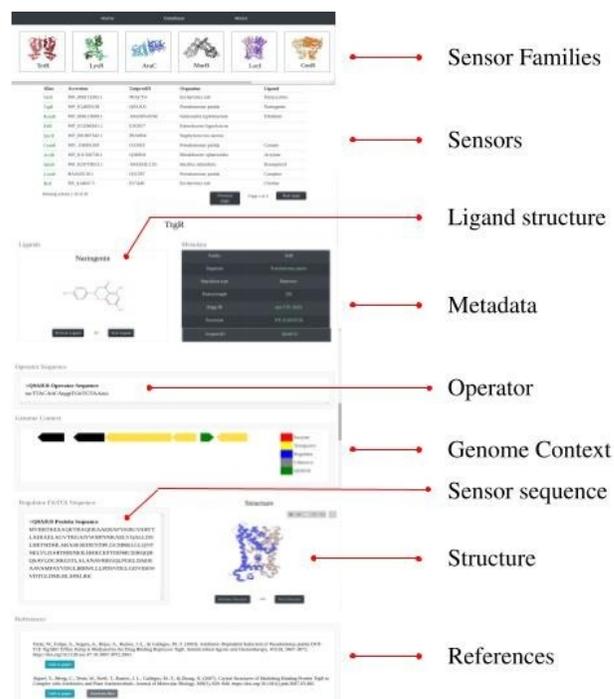


Figure 2: Layout of a typical biosensor entry.

Web application architecture

All sensor data is stored in a SQLite relational database and accessed by a Flask server via SQLAlchemy. React is used on the front-end to handle user requests, interface with the Flask back-end, and dynamically update the information displayed. The bootstrap framework is used for styling and external plugins are used for chemical structure display (smilesDrawer) and protein structure display (LiteMol). The application is deployed using Docker containers on an AWS server.

This database can be accessed at <https://gbiosensors.com>

3 ACKNOWLEDGEMENTS

The authors are affiliated with UT Austin, but have not received any university or external support for this project.

REFERENCES

- [1] ET AL., M. R. High-throughput screening for high-efficiency small-molecule biosynthesis. *Metab Eng* 63, 1 (2020), 102–125.
- [2] JIE ZHANG, MICHAEL K JENSEN, J. D. K. Development of biosensors and their application in metabolic engineering. *Curr Opin Chem Biol* 28, 10 (2015), 1–8.
- [3] PIERRE KUGLER, DEBORAH FRÖHLICH, V. F. W. Development of a biosensor for crotonobetaine-coa ligase screening based on the elucidation of escherichia coli carnitine metabolism. *ACS Synth Biol* 9, 9 (2020), 2460–2471.
- [4] YOUNGCHANG KIM, GRAZYNA JOACHIMIAK, L. B. G. B. A. J. How aromatic compounds block dna binding of hcr catabolite regulator. *J Biol Chem* 291, 25 (2016), 13243–56.

Improving ensemble model predictions in a general biosynthesis experiment modeling tool

Bret Peterson

bretpeterson@lbl.gov
Lawrence Berkeley Laboratory
Berkeley, CA

Tijana Radivojevic

tradivojevic@lbl.gov
Lawrence Berkeley Laboratory
Berkeley, CA

Hector Garcia Martin

hgmartin@lbl.gov
Lawrence Berkeley Laboratory
Berkeley, CA

1 INTRODUCTION

Synthetic biology holds promise for engineering solutions for a large variety of world-wide problems including therapeutics for novel diseases [10], carbon neutral fuels [4], and new materials for reducing environmental impact [2, 8]. Given the vast parameter space of possible but fruitless experimental manipulations, the pace of discovery in synthetic biology is tied to our ability to design experiments with a high likelihood of success. Complex, nonlinear interactions between sources, intermediates, products and control systems throughout the cellular processes make it difficult to optimize the production through rational design. Improving production is instead heavily reliant on empirical results.

Machine learning can help guide empirically based experimental design [1, 5, 11]. However, unlike in other fields, such as text or image processing which can have hundreds of millions of training examples, training examples in synthetic biology are often restricted to what highly skilled researchers can generate through often expensive and time-consuming experiments.

Our group has been developing the Automated Recommendation Tool (ART) for making predictions of experimental results and using these predictions to recommend new experiments based on the smaller training data sets representative of many synthetic biology development efforts. ART has already proven to be useful for a variety of synthetic biology use cases, including production of renewable biofuels, hoppy flavored beer, fatty acids, and tryptophan [9, 11]. The predictions are generated based on an ensemble of models, and the probability distributions of weights on the underlying models are used to characterize uncertainty in the predictions from the final, ensemble model. This uncertainty characterization is a key component of ART and is important in multiple phases of optimization as it indicates areas of inadequate exploration in the early phases, and areas of likely success for goal-driven exploitation in the later phases of optimization.

Because the nature of the contribution of the multiple models in the ensemble to the predictions is critical to both accuracy and the indicated uncertainty, we are exploring the distribution of contributions of multiple models to the end results. In this paper, we characterize preliminary results

on the impact of increasing the number of models in the ensemble that are based on optimized model selection and hyperparameter/workflow tuning (TPOT [7]) in order to improve accuracy of the ensemble predictions and increase the chances that multiple models can substantially contribute to those predictions.

2 METHODS

For this preliminary study, we used an explicit function so that there was a ground truth upon which to judge performance directly. We used a function that had been previously used as a challenging function for models to learn [3, 9]. The function has a large number of local minima across the feature space. For a d -dimensional vector x , the function is

$$F(x) = \sum_{i=1}^d \sqrt{x_i} \sin(x_i). \quad (1)$$

A training set of 1024 samples was generated using Latin Hypercube [6] to select feature values from a domain spanning (0, 12) for each feature. Seven fixed regression model types with default hyperparameters from the well-known Scikit-learn library were trained as well as a variable number of models (1, 2, or 4) based on a TPOT optimized workflow (Figure 1) that selects a model type and associated hyperparameters according to performance on the training set. Weights for the results of each model were trained to create a Bayesian ensemble model, and the resulting ensemble model was used to make predictions for the training and test sets. The test set was generated similarly to the training set with 1024 Latin Hypercube samples. Accuracy was assessed based on the coefficient of determination for the mean values of predictions relative to the calculated (true) values for both the training and test sets. Experiments were repeated for a different number of features in the feature space (4, 8, or 16), which provided an increasing level of difficulty for accurately representing the function, both because of increased complexity of the function and reduced sampling relative to the number of features. The experiments were repeated 4 times each to gauge variation due to differences in the sampling of the feature space for the training and test sets.

3 RESULTS

The results for prediction accuracy are shown in the bar graphs of Figure 2. A value of 1.0 indicates perfect prediction, and a value of 0.0 indicates that the variance in the errors is of equal size as the variance in the predicted values. When there were only 4 features, the models are able to accurately capture the underlying function and there is little difference between the training and test sets' results regardless of the number of TPOT-based models that were added. As the number of features was increased to 8 and then 16, the function became much more difficult to learn. Consequently, the prediction accuracy was reduced for both the training and test sets, and the difference in prediction accuracy between the training and test sets became more marked. Under these conditions, additional TPOT models improved the accuracy. In the case of 16 features, the means of the coefficients of determination for 1, 2, 4 TPOT-based models was 0.812, 0.875, 0.890 for the test data set and 0.942, 0.970, 0.975 for the training data set respectively.

4 DISCUSSION

When using ART, we were not surprised that the model generated from the TPOT workflow generation would sometimes be very heavily weighted with respect to the other trained models since it involves optimization of model selection and its hyperparameters before training. This heavy reliance on the TPOT-based model is not an issue when that model accurately captures the underlying function in all areas of interest of the feature space. In such cases, the ensemble approach is not expected to provide substantial benefit.

For more challenging underlying functions, an ensemble model can improve performance by providing a weighted average of model outputs. This weighted average can reduce random errors across the models, as well as minimize an outlier in a single model that contradicts the other models in a particular region of the feature space.

Towards this goal, we introduced multiple models generated from different TPOT model selection and hyperparameter/workflow tuning runs using different random seeds. The model selection and hyperparameter tuning is sufficiently sensitive to initial conditions that we were able to inspect the generated workflows and confirm that these workflows were indeed significantly different, often with different underlying types of regression models. The weights assigned to the multiple TPOT-based models were more evenly distributed when the underlying function was made more difficult for a single model to represent (pie charts in Figure 2 indicate these distributions). The prediction accuracy also improved relative to the standard ART use case of a single TPOT-based model.

Since ART's performance on this function and on actual synthetic biology predictions have already been demonstrated [9, 11], we suggest that improving performance in ART on this function indicates likely improvements on challenging biosynthesis experimental predictions as well.

The demonstrated improvement in accuracy and suspected better characterization of uncertainty should result in improved experimental recommendations. Such improvements will more quickly lead to successful experimental outcomes, reducing experimental time and costs, and ultimately resulting in viable synthetic biology products.

5 FUTURE WORK

We plan to more fully characterize the impact of more even distribution of weights in the ensemble model on the calculated uncertainty and how that uncertainty reflects true values for different classes of functions. Also, given the assumption of independence between models in the ensemble model, we plan to explore more explicit sampling without replacement in the TPOT model selection optimization which would ensure that the base model was always different for each recommended workflow. We will then use this approach in some of our biosynthesis projects and will make it generally available in subsequent versions of ART if it proves valuable as these results suggest.

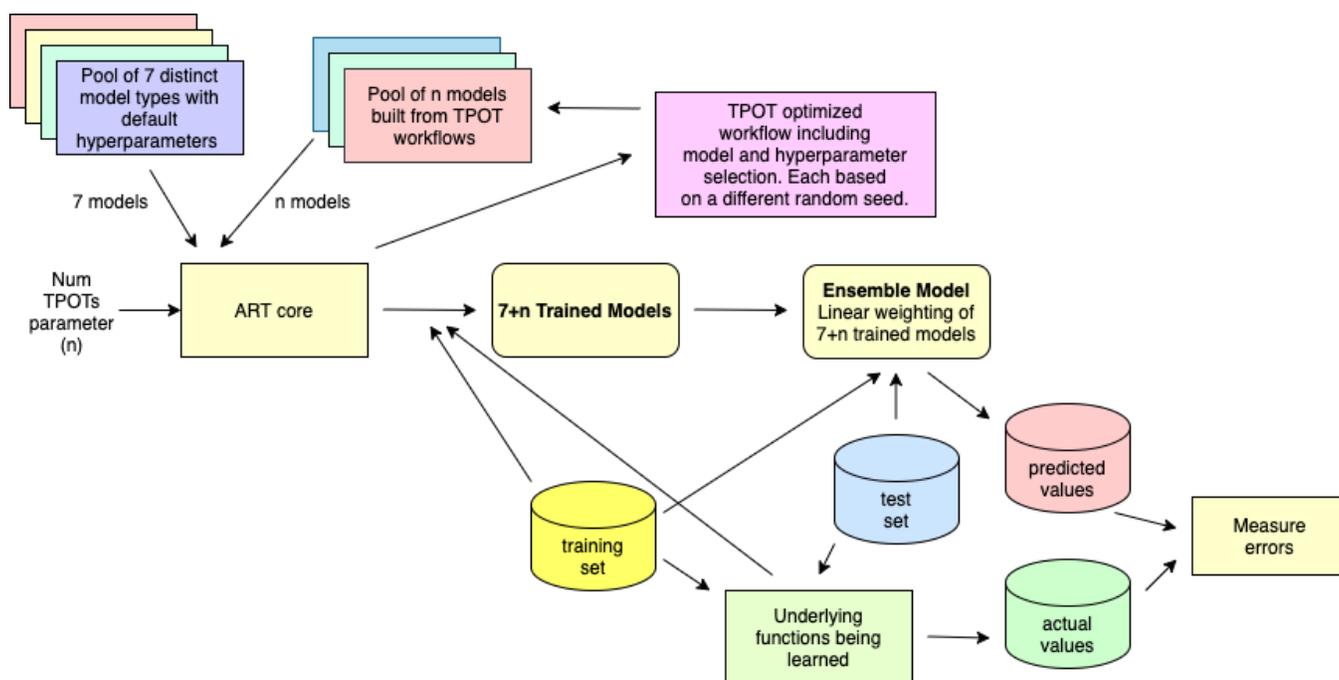


Figure 1: Overview of selection of models for a given run to create an ensemble model. The models and ensemble weights were generated based on a training set for a feature space of variable size, and the ensemble model predictions were assessed based on the test set.

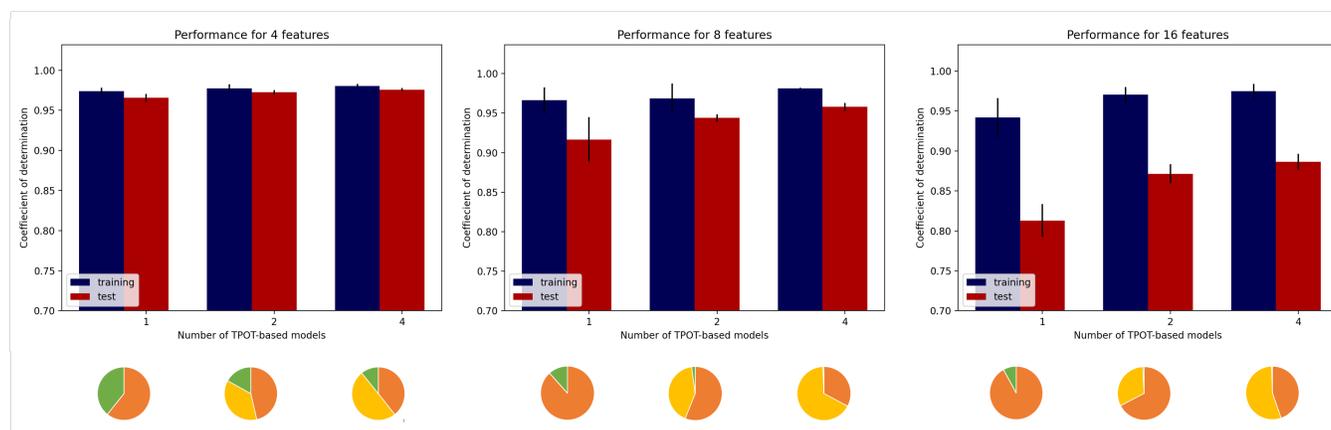


Figure 2: Prediction accuracy improves more markedly with more TPOT-based models as the number of features increases. The bar graphs show the prediction accuracy based on the coefficient of determination for ensembles including 1, 2, and 4 TPOT-based models. The blue bars indicate performance on the training data sets, and the red bars indicate performance on the test data sets. The three graphs are for 4, 8, and 16 features, representing increasingly difficult functions to learn. The pie graphs under the pairs of bars for each ensemble model indicate the distribution of the weighting of the individual models within the ensemble. Orange indicates the proportion for the most heavily weighted TPOT-based model, and green indicates the proportion for the 7 Scikit-learn regression models. For the right two columns in each graph, the yellow indicates the proportion for the additional TPOT-based model (middle) or the combined weight of the other 3 TPOT-based models (rightmost). More even distribution of weights across multiple models within the ensemble model correlates with improved prediction accuracy.

REFERENCES

- [1] CARBONELL, P., RADIVOJEVIC, T., AND GARCÍA MARTÍN, H. Opportunities at the intersection of synthetic biology, machine learning, and automation. *ACS Synth Biol* 8, 7 (jul 2019), 1474–1477.
- [2] JAIN, R., AND TIWARI, A. Biosynthesis of planet friendly bioplastics using renewable carbon source. *J Environ Health Sci Eng* 13 (feb 2015), 11.
- [3] JAMIL, M., AND YANG, X. S. A literature survey of benchmark functions for global optimisation problems. *IJMMNO* 4, 2 (2013), 150.
- [4] KEASLING, J., GARCIA MARTIN, H., LEE, T. S., MUKHOPADHYAY, A., SINGER, S. W., AND SUNDBLUM, E. Microbial production of advanced biofuels. *Nat Rev Microbiol* (jun 2021).
- [5] LAWSON, C. E., MARTÍ, J. M., RADIVOJEVIC, T., JONNALAGADDA, S. V. R., GENTZ, R., HILLSON, N. J., PEISERT, S., KIM, J., SIMMONS, B. A., PETZOLD, C. J., SINGER, S. W., MUKHOPADHYAY, A., TANJORE, D., DUNN, J. G., AND GARCIA MARTIN, H. Machine learning for metabolic engineering: A review. *Metab Eng* 63 (jan 2021), 34–60.
- [6] MCKAY, M. D., BECKMAN, R. J., AND CONOVER, W. J. A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics* 21, 2 (may 1979), 239.
- [7] OLSON, R. S., AND MOORE, J. H. TPOT: A tree-based pipeline optimization tool for automating machine learning. In *Automated machine learning: methods, systems, challenges*, F. Hutter, L. Kotthoff, and J. Vanschoren, Eds., The springer series on challenges in machine learning. Springer International Publishing, Cham, 2019, pp. 151–160.
- [8] PODDAR, H., BREITLING, R., AND TAKANO, E. Towards engineering and production of artificial spider silk using tools of synthetic biology. *Eng. biol.* 4, 1 (mar 2020), 1–6.
- [9] RADIVOJEVIĆ, T., COSTELLO, Z., WORKMAN, K., AND GARCIA MARTIN, H. A machine learning automated recommendation tool for synthetic biology. *Nat Commun* 11, 1 (sep 2020), 4879.
- [10] WEBER, W., AND FUSSENEGGER, M. The impact of synthetic biology on drug discovery. *Drug Discov Today* 14, 19-20 (oct 2009), 956–963.
- [11] ZHANG, J., PETERSEN, S. D., RADIVOJEVIC, T., RAMIREZ, A., PÉREZ-MANRÍQUEZ, A., ABELIUK, E., SÁNCHEZ, B. J., COSTELLO, Z., CHEN, Y., FERRO, M. J., MARTIN, H. G., NIELSEN, J., KEASLING, J. D., AND JENSEN, M. K. Combining mechanistic and machine learning models for predictive engineering and optimization of tryptophan metabolism. *Nat Commun* 11, 1 (sep 2020), 4880.

iBioSim Server: a Tool for Improving the Workflow for Genetic Design and Modeling

Thomas Stoughton¹, Lukas Buecherl¹, Payton J Thomas², Pedro Fontanarrosa², Chris J. Myers¹

¹University of Colorado Boulder, ²University of Utah
chris.myers@colorado.edu

1 INTRODUCTION

Synthetic biology is the field of research that deals with the application of engineering principles to design biological systems to perform user desired functions [2]. Research into this field has many applications including the development of bio-fuels [8], internal drug-delivery systems [10], and bio-sensors [12], all of which rely on genetic circuits to work. Genetic circuits use the proteins and cellular pathways of living cells to execute their functions.

Research in synthetic biology—particularly developing genetic circuits—usually follows a *Design, Build, Test, Learn* (DBTL) cycle to progress from a desired application to a physical build. In order to aid in the design of genetic circuits, computational models are used to simulate them before attempts are made to build them *in vivo*. This saves time for researchers and can help to predict malfunctions in the circuit before they are built.

Genetic design automation (GDA) tools to support the DBTL cycle include SynBioHub [7], SBOLCanvas [11], and iBioSim [14]. SynBioHub is an online repository for sharing genetic parts and other design information. SBOLCanvas is a web application that can be used to create genetic circuits using parts stored in SynBioHub (<https://sbolcanvas.org/>). Finally, the iBioSim application is a software tool for the design, modeling, and analysis of genetic circuits.

A set of standard description languages are used to link these tools together. These standards allow for coherent communication between the tools and for the reproducibility of results that is necessary in synthetic biology. The languages used by the aforementioned software are the *Synthetic Biology Open Language* (SBOL) [4, 6, 9], the *Systems Biology Markup Language* (SBML) [5], and the *Simulation Experiment Design Markup Language* (SED-ML) [13]. SBOL, SBML, and SED-ML are all XML-based languages, each of which has a different purpose for the design and simulation of genetic circuits. An SBOL file encodes genetic parts and their relationships within a circuit, SBML files encode computational models of genetic circuits, and SED-ML files encode the simulation specifications for the purpose of experiment repeatability. In addition to these languages, a COMBINE Archive [1] is used to package SBOL, SBML, and SED-ML files together so that a simulation study can be replicated easily.

With the development of these software tools and standards, a general workflow for designing, modeling, and simulating a genetic circuit was established. The current workflow for using these software tools is specified in Figure 1(a). The process of designing and testing genetic circuits starts with the retrieval of genetic parts from SynBioHub through SBOLCanvas. Here, the initial design for a genetic circuit can be created. Then, the VirtualParts API [3] enriches the DNA-level design with additional proteins and small molecules that are known to interact with the genetic parts. Next, the computational model is created in iBioSim, where the SBOL file automatically gets converted into an SBML file. Then, the simulation conditions are decided and the simulation is run. Finally, the results of the simulation are rendered and then sent back to SynBioHub for sharing.

This workflow is a complex, manual process of importing and exporting files when switching between applications, so automating aspects of this process is becoming a priority. In addition, iBioSim’s graphical user interface is outdated and not easy to follow for someone not already familiar with it.

These problems could be solved with some level of design automation. One of the biggest steps forward would be to keep researchers on the newer SBOLCanvas website throughout the entire workflow. This removes the need to download and configure an application, keeps the entire process on a web-based system, and only uses the more intuitive user interface of SBOLCanvas.

2 RESULTS

In order to take advantage of the simulation capabilities of the iBioSim application, this work aims to push these functional aspects to the back-end of the SBOLCanvas website. This can be accomplished using an API that communicates with SBOLCanvas and is capable of running iBioSim on a server. A proposed workflow for using an API is described in Figure 1(b). The long-term goal for this API is for it to be able to receive an http request from SBOLCanvas that contains an SBOL file and arguments for completing the Enrich, Model, and Simulate steps, and then execute those steps automatically on the server. The results from these steps would then be packaged into a COMBINE Archive, and then sent back to SBOLCanvas. The entire API would be running on a Docker container to make it easier to update externally. A Docker

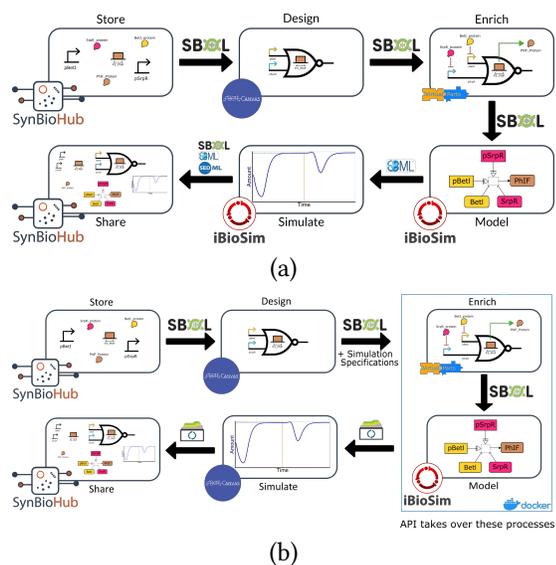


Figure 1: (a) Diagram of the current workflow for designing and simulating a genetic circuit. This workflow involves three different applications: the Design step happens in SBOLCanvas, the Enrichment, Model, and Simulate steps happen within iBioSim, and the Share step happens in SynBioHub. At each change of application, the relevant SBOL, SBML, or SED-ML files must be exported from the previous application and then imported into the next application. (b) Diagram of the new proposed workflow. This workflow only requires the user to work on the SBOLCanvas application; the Enrichment, Model, and the Simulate step are taken over by the API provided by the iBioSim server containerized using Docker, and the simulation is rendered back on SBOLCanvas. This workflow also makes use of the COMBINE Archive to package the SBOL, SBML, and SED-ML files together for ease-of-use and for reproducibility.

container is a virtual environment that stores an application for standalone execution. This approach eliminates the need for researchers to download an application, while keeping all of the same capabilities. The ability to keep the entire user-side process on a single web-application would undoubtedly provide for a more seamless workflow when compared to the current workflow.

The iBioSim server application is a Dockerized API implemented using Python that is capable of receiving an http request with a COMBINE Archive from a previously completed simulation experiment, reproducing the results, and sending the results back as a response to the http request. The API accomplishes this by running the executable iBioSim analysis file which can complete a simulation from the command line. The API then collects the results generated from the simulation and sends them back, after which the unnecessary files are deleted to prepare for another http request.

For the purposes of testing, the API only received COMBINE Archives from previously generated simulation experiments as an input so that the results from the API could be compared to known results of the simulation. In addition, since SBOLCanvas has yet to be updated to be able to send the correct http requests, the tests were sent using the API development tool, Postman.

3 DISCUSSION

With this work, a containerized API was created that is capable of automating the simulation of genetic circuits designed in SBOLCanvas. The work on this API shows promise for realizing the new, streamlined workflow for virtually developing and testing genetic circuits. Web-based applications have many advantages for both users and developers over standard, downloadable applications. Some of these advantages include accessibility, ease of maintenance, and the elimination of installation issues. There are a few drawbacks of a completely web-based approach, such as the need for internet access and the unavailability of iBioSim's code, but we believe these trade-offs are favorable for the average user. This approach would allow researchers to have a single, entirely online work-space for designing and simulating genetic circuits, something which has only previously been available in downloadable application packages. The idea of dedicating a server to running simulations and reusing iBioSim's functional parts saves a lot of effort compared to a full overhaul of iBioSim onto a web-application. With this approach, no functionality in iBioSim other than the GUI is lost as the program was designed to be fully functional through the command line, though the dynamic modeling of circuits has yet to be implemented. Updating the API to stay in sync with changes to iBioSim is trivial due to the nature of the Docker container, which uses the newest published version of iBioSim upon building. Only the code for the API would need to be updated, and the Docker container would simply be restarted.

However, there are some more significant steps to be made before this API is ready for public use. As previously mentioned, the API is currently only able to accept COMBINE Archives as an input file. In the future, it should handle single SBOL files as input so that first-time simulations can be run. For this, the iBioSim server will be expanded to include the SBOL-to-SBML conversion aspect from iBioSim. In addition, SBOLCanvas will need to be updated to send http requests to the API with the proper information to allow for seamless front-end simulations. The logistics of setting up a server that can run complex simulations in a timely manner are also being considered for the future. The code for this API is open source, available on [Github](#), and released with the Apache-2.0 License.

Acknowledgements

The authors of this work are supported by DARPA Grant FA8750-17-C-0229, National Science Foundation Grant No. 1856740, and the University of Colorado, Boulder Engineering Excellence Fund. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the funding agencies.

REFERENCES

- [1] BERGMANN, F. T., ADAMS, R., MOODIE, S., COOPER, J., GLONT, M., GOLEBIEWSKI, M., HUCKA, M., LAIBE, C., MILLER, A. K., NICKERSON, D. P., ET AL. Combine archive and omex format: one file to share all information to reproduce a modeling project. *BMC bioinformatics* 15, 1 (2014), 1–9.
- [2] CHENG, A. A., AND LU, T. K. Synthetic biology: An emerging engineering discipline. *Annual Review of Biomedical Engineering* 14, 1 (2012), 155–178. PMID: 22577777.
- [3] COOLING, M. T., ROUILLY, V., MISIRLI, G., LAWSON, J., YU, T., HALLINAN, J., AND WIPAT, A. Standard virtual biological parts: a repository of modular modeling components for synthetic biology. *Bioinformatics* 26, 7 (2010), 925–931.
- [4] GALDZICKI, M., CLANCY, K. P., OBERORTNER, E., POCOCK, M., QUINN, J. Y., RODRIGUEZ, C. A., ROEHNER, N., WILSON, M. L., ADAM, L., ANDERSON, J. C., ET AL. The synthetic biology open language (SBOL) provides a community standard for communicating designs in synthetic biology. *Nature biotechnology* 32, 6 (2014), 545–550.
- [5] KEATING, S. M., WALTEMATH, D., KÖNIG, M., ZHANG, F., DRÄGER, A., CHAOUIYA, C., BERGMANN, F. T., FINNEY, A., GILLESPIE, C. S., HELIKAR, T., ET AL. SBML level 3: an extensible format for the exchange and reuse of biological models. *Molecular systems biology* 16, 8 (2020), e9110.
- [6] McLAUGHLIN, J. A., BEAL, J., MISIRLI, G., GRÜNBERG, R., BARTLEY, B. A., SCOTT-BROWN, J., VAIDYANATHAN, P., FONTANARROSA, P., OBERORTNER, E., WIPAT, A., ET AL. The synthetic biology open language (SBOL) version 3: simplified data exchange for bioengineering. *Frontiers in Bioengineering and Biotechnology* 8 (2020), 1009.
- [7] McLAUGHLIN, J. A., MYERS, C. J., ZUNDEL, Z., MISIRLI, G., ZHANG, M., OFITERU, I. D., GONI-MORENO, A., AND WIPAT, A. SynBioHub: a standards-enabled design repository for synthetic biology. *ACS synthetic biology* 7, 2 (2018), 682–688.
- [8] OLSON, D. G., MCBRIDE, J. E., JOE SHAW, A., AND LYND, L. R. Recent progress in consolidated bioprocessing. *Current Opinion in Biotechnology* 23, 3 (2012), 396–405. Energy biotechnology • Environmental biotechnology.
- [9] ROEHNER, N., BEAL, J., CLANCY, K., BARTLEY, B., MISIRLI, G., GRUNBERG, R., OBERORTNER, E., POCOCK, M., BISSELL, M., MADSEN, C., ET AL. Sharing structure and function in biological design with SBOL 2.0. *ACS synthetic biology* 5, 6 (2016), 498–506.
- [10] SHI, P., GUSTAFSON, J. A., AND MACKAY, J. A. Genetically engineered nanocarriers for drug delivery. *International journal of nanomedicine* 9 (2014), 1617.
- [11] TERRY, L., EARL, J., THAYER, S., BRIDGE, S., AND MYERS, C. J. Sbolcanvas: A visual editor for genetic designs. *ACS Synthetic Biology* (2021).
- [12] VIGNESHVAR, S., SUDHAKUMARI, C., SENTHILKUMARAN, B., AND PRAKASH, H. Recent advances in biosensor technology for potential applications—an overview. *Frontiers in bioengineering and biotechnology* 4 (2016), 11.
- [13] WALTEMATH, D., ADAMS, R., BERGMANN, F. T., HUCKA, M., KOLPAKOV, F., MILLER, A. K., MORARU, I. I., NICKERSON, D., SAHLE, S., SNOEP, J. L., ET AL. Reproducible computational biology experiments with SED-ML—the simulation experiment description markup language. *BMC systems biology* 5, 1 (2011), 1–10.
- [14] WATANABE, L., NGUYEN, T., ZHANG, M., ZUNDEL, Z., ZHANG, Z., MADSEN, C., ROEHNER, N., AND MYERS, C. ibiosim 3: a tool for model-based genetic circuit design. *ACS synthetic biology* 8, 7 (2018), 1560–1563.

SynBioHub2 - Providing an Intuitive and Maintainable Genetic Design Repository

Benjamin Hatch¹, Jeanet Mante², Eric Yu¹, Chris J. Myers²

¹University Of Utah, ²University of Colorado Boulder
chris.myers@colorado.edu

1 INTRODUCTION

One of the primary goals of the field of synthetic biology is to make genetic designs easier to develop. In order to achieve this goal, genetic design repositories, such as JBEI-ICE [1] and the *International Genetically Engineered Machine* (iGEM) Registry of Standard Biological Parts (<http://parts.igem.org>), have been developed to facilitate storing and sharing genetic designs. These repositories are useful for DNA sequence sharing, but they provide little information about genetic designs at an abstract level, such as information about proteins, interactions, and design versioning. In order to address the need for supplying information about synthetic designs at the abstract (as well as the sequence) level, SynBioHub [2] was developed. SynBioHub is an open-source, web-based design repository that facilitates storing, searching, and sharing of genetic designs. It utilizes the Synthetic Biology Open Language (SBOL) [3] data format to store both sequence and abstract metadata of genetic designs.

Although SynBioHub is already utilized by many synthetic biologists and organizations, further development is necessary to meet the needs of large-scale synthetic biology projects. This development has proved to be challenging, as SynBioHub's code base is difficult to maintain due to lack of software modularity and outdated technologies. To solve this issue, SynBioHub's back and front-end architecture is being redesigned to provide a more maintainable and intuitive genetic design repository that will be resilient to new functionality development and future SBOL version releases. The development of the redesigned back-end and front-end is happening in parallel; development of the front-end is almost complete.

This paper presents SynBioHub2's redesigned front-end, which will be soon be released as SynBioHub2. This version of SynBioHub uses the same back-end as the original SynBioHub while providing a more intuitive and maintainable user interface than that of the original SynBioHub.

2 RESULTS

By redesigning SynBioHub's front-end, SynBioHub2 provides a more intuitive user experience for synthetic biologists. This is illustrated through SynBioHub2's new submit, search, and viewing interfaces.

Submitting Designs

Given that SynBioHub's primary purpose is to facilitate the storage of genetic parts, it is important that SynBioHub provides an intuitive submit workflow for its users. SynBioHub2's submit workflow has been redesigned to make submission more intuitive. The current SynBioHub prevents users from submitting more than one design file at a time. If users wish to upload more than one design file, they must either upload each file individually or zip their design files. When submitting a zip file, the submission will fail if any of the design files cannot be processed. SynBioHub2's submit interface allows multiple files to be submitted at once. By utilizing React, the state of each file's submission is dynamically displayed. Users can also continue to use other parts of the application (such as browsing design collections) while waiting for their submission to complete. If any of the design files fail to upload, a specific error message is displayed that explains why the file failed, and users are given the option to submit the file as a collection attachment. SynBioHub2's submit interface has already been developed; a screenshot of the interface is shown in Figure 1.

Searching for Designs

One of the main purposes of SynBioHub is to enable users to efficiently search for genetic designs. SynBioHub2's search interface has been improved as follows. First, SynBioHub2 provides feedback when searching for a genetic design. Its interface uses the React JavaScript framework to dynamically display animations based on the state of a search's execution, such as a loading icon when the search is being processed and relevant error messages when a search fails. Second, SynBioHub2 displays search results in a format that resembles a spreadsheet. This allows users to view more part results at once, as well as makes the search experience more intuitive for synthetic biologists, who commonly utilize spreadsheets to store information about their genetic designs. Lastly, SynBioHub2 enables users to quickly perform frequently used actions on search results. These actions include downloading design files, creating a design collection, and storing search results for later reference in a virtual "basket". SynBioHub2's search interface has already been developed; a screenshot of the interface can be found in Figure 2.

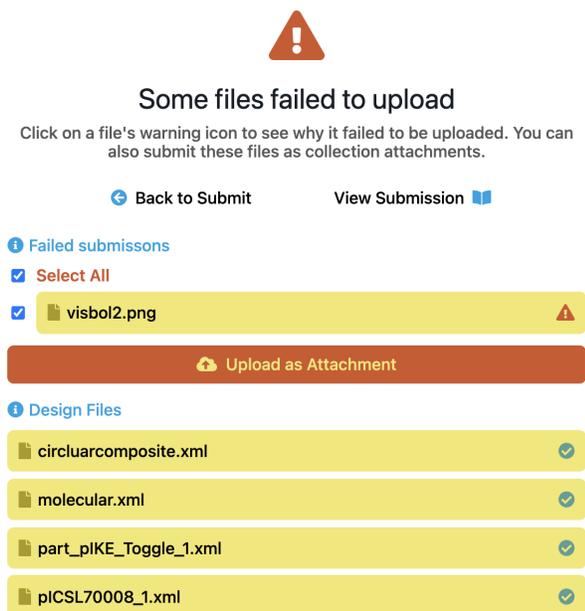


Figure 1: SynBioHub2’s submit interface. Users can view the status of each submitted design file and are given the option to submit files as attachments when they fail to upload. In this case, the user attempted to submit an image as a design file, which is not supported by SynBioHub. The user is notified that this submission failed, and is prompted to submit the file as a collection attachment.

Viewing Designs

The main issue that prompted the redesign of SynBioHub was an outdated interface for viewing genetic parts. The current SynBioHub’s front-end parses the entire file representing the design that the user wishes to view before serving an overview and renderings of the design. This has proven to be detrimental to the user experience when users view larger design files, as the viewing interface often takes a long time to load. Additionally, adding any native code to the viewing interface often requires further parsing of design files, which further increases page load times.

SynBioHub2’s viewing interface will tackle the viewing of genetic designs by using a dynamic rendering pipeline. Rather than parse entire design files at once, SynBioHub2 will first make a request to its back-end to fetch relevant metadata about the design to display to the user. This metadata will be used to populate the viewing page, which will be served to the corresponding user immediately following this metadata population. After the viewing page is served to the user, aspects of the viewing page that require further parsing of the page’s corresponding design file will be asynchronously

fetches and rendered. This asynchronous rendering pipeline will significantly reduce the time it takes for users to view genetic designs. Additionally, SynBioHub2’s viewing interface will give users control over which viewing elements are visible, as well as the order in which the elements are displayed through an interactive page overview pane. This viewing functionality is still under active development. A design mock-up for SynBioHub2’s viewing interface can be found in Figure 3.

3 METHODS

SynBioHub2 achieves greater maintainability and scalability by implementing a separation of concerns. It strictly divides its front-end code from its back-end. This allows SynBioHub2 to subscribe to the original SynBioHub’s back-end until the new back-end’s development is complete. SynBioHub2’s front-end refactors UI patterns into reusable React components that can be modified quickly, enabling new features and changes to be implemented relatively seamlessly. It also features more extensive code documentation, which will enable developers to quickly contribute to its open-source code base without introducing software bugs in the application.

4 DISCUSSION

Although SynBioHub2 makes a significant number of improvements to the original SynBioHub, it will not be the final version of SynBioHub. Rather, it is a stepping stone for SynBioHub3, which has a number of design goals, which are listed below.

- Community involvement throughout development.
- A more intuitive front-end for biologists.
- A faster, more flexible user experience.
- Improved integration of curation.
- Support for SBOL3 and beyond.
- Preservation of the existing back-end and plugin API.

SynBioHub2 already begins to fulfill some of these design goals; it provides a more intuitive and interactive front-end for biologists, and community involvement throughout the development process has been encouraged. However, design goals such as the ability to use different triplestore databases are out of its scope, as SynBioHub2 focuses mostly on redesigning of the front-end. These design goals will be achieved once SynBioHub2 transfers its subscription from SynBioHub’s original back-end to its redesigned back-end, which is being developed in parallel to SynBioHub2.

SynBioHub2 is under active development. Although a significant portion of the application has been developed, such as submitting and searching for parts, some functionality still needs to be implemented. Such functionality includes design viewing, sharing of designs, and support for design viewing plugins. Users can experiment with SynBioHub2’s

The screenshot shows the SynBioHub2 search interface. At the top, there is a search bar with the query 'gfp'. Below the search bar, there are navigation links: 'Standard Search', 'Root Collections', 'Sequence Search', 'Advanced Search', and 'SPARQL'. There are also buttons for '+ Add to Basket' and 'Download'. The search results are displayed in a table with 5 columns: Name, Display ID, Description, Type, and Privacy. The first result is 'GFP' with Display ID 'BBa_E0040' and Description 'green fluorescent protein derived from jellyfish Aequorea victoria wild-type GFP (SwissProt: P42212)'. There are 15 results in total, with the first one selected.

<input type="checkbox"/>	Name	Display ID	Description	Type	Privacy
<input checked="" type="checkbox"/>	GFP	BBa_E0040	green fluorescent protein derived from jellyfish Aequorea victoria wild-type GFP (SwissProt: P42212)	Component	🔒
<input type="checkbox"/>	GFP report	BBa_E0240	GFP generator	Component	🔒
<input type="checkbox"/>	BBa_J04450	BBa_J04450	RFP Coding Device	Component	🔒
<input type="checkbox"/>	YFP	BBa_E0032	enhanced yellow fluorescent protein derived from A. victoria GFP	Component	🔒
<input type="checkbox"/>	eyfp	BBa_E0030	enhanced yellow fluorescent protein derived from A. victoria GFP	Component	🔒
<input type="checkbox"/>	GFP genera	BBa_E0840	GFP generator	Component	🔒
<input type="checkbox"/>	ecfp	BBa_E0020	engineered cyan fluorescent protein derived from A. victoria GFP	Component	🔒
<input type="checkbox"/>	ECFP	BBa_E0022	enhanced cyan fluorescent protein derived from A. victoria GFP	Component	🔒
<input type="checkbox"/>	mCherry	BBa_J06504	monomeric RFP optimized for bacteria	Component	🔒
<input type="checkbox"/>	GFP	BBa_K145015	GFP with LVA tag	Component	🔒
<input type="checkbox"/>	BBa_K208000	BBa_K208000	GFP	Component	🔒
<input type="checkbox"/>	GFP1	BBa_I715019	Amino Half of GFP (aka GFP1)	Component	🔒
<input type="checkbox"/>	GFP2	BBa_I715020	Carboxyl Half of GFP (aka GFP2)	Component	🔒
<input type="checkbox"/>	RBS-RFP	BBa_K093005	RFP with RBS	Component	🔒

Figure 2: SynBioHub2’s search interface. Users can check search results and perform common actions, such as downloading corresponding parts.

development version at: <https://dev2.synbiohub.org/>. SynBioHub2’s open-source code base can be found on GitHub at: <https://github.com/SynBioHub/synbiohub3/>.

5 ACKNOWLEDGEMENTS

The authors of this work are supported by DARPA FA8750-17-C-0229, the National Science Foundation Grant No. 1939892, and a Dean’s Graduate Assistantship at the University of Colorado Boulder. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the funding agencies.

REFERENCES

- [1] HAM, T., DMYTRIV, Z., PLAHAR, H., CHEN, J., HILLSON, N., AND KEASLING, J. Design, implementation and practice of jbei-ice: an open source biological part registry platform and tools. *Nucleic acids research* 40, 18 (2012), e141–e141.
- [2] McLAUGHLIN, J. A., MYERS, C. J., ZUNDEL, Z., MISIRLI, G., ZHANG, M., OFITERU, I. D., GOÑI-MORENO, A., AND WIPAT, A. SynBioHub: A standards-enabled design repository for synthetic biology. *ACS Synthetic Biology* 7, 2 (2018), 682–688. PMID: 29316788.
- [3] ROEHNER, N., BEAL, J., CLANCY, K., BARTLEY, B., MISIRLI, G., GRÜNBERG, R., OBERORTNER, E., POCOCK, M., BISSELL, M., MADSEN, C., NGUYEN, T., ZHANG, M., ZHANG, Z., ZUNDEL, Z., DENSMORE, D., GENNARI, J. H., WIPAT, A., SAURO, H. M., AND MYERS, C. J. Sharing structure and function in biological design with SBOL 2.0. *ACS Synthetic Biology* 5, 6 (2016), 498–506. PMID: 27111421.

GFP Report (BBa_E0240)
Component, Version 1

BBa_E0240
Terminator, Circular
http://parts.igem.org/Part:BBa_E0240
Jennifer Braff
2004-10-17 11:00:00
<https://synbiohub.org/public/igem/igem2/sbol/>
Synthetic (32630)
E. Coli (562)
1560 Uses, 75 twins, 16 similar

Page Sections

- VisBOL
- Description
- Sequence Visualization
- Other Plugin
- Further Properties
- Sequence
- Citations and References

VisBOL

BBa_B0032 GFP BBa_B0010 BBa_B0012

Description

GFP reporter used for the testing of different promoter strengths. This is meant to be filled in with a general description that might contain information similar to the wiki

Design Notes
This sequence was adapted from the source by codon optimisation. The initial design was for S. pombe and so the reduction of U codons was carried out.

Composition Type
Uses the MoClo Assembly Method/ Was Synthesised, etc

Sequence Visualization

Figure 3: Mock-up for SynBioHub2’s design view. Note the overview pane on the left and the ability to check/uncheck which sections of the page are visible.

An Investigative Platform Comprising Cell-Free Transcription- Translation and Electron Microscopy for Studying Bacteriophages

Joseph P. Wheatley^{1,2}, Sahan B. W. Liyanagedera^{1,2}, Ian Hands-Portman², Antonia P. Sagona², Vishwesh Kulkarni^{1*}

¹University of Warwick, School of Engineering, ²University of Warwick, School of Life Sciences
{Joseph.Wheatley.1,S.Liyanagedera,I.J.Portman,A.Sagona,V.Kulkarni}@warwick.ac.uk

1 INTRODUCTION

Understanding the structure and propensity of bacteriophages (phages) can lead to key insights on how these highly abundant viruses can be harnessed and exploited for use in therapeutics and diagnostics. Cell-free transcription-translation (TXTL) systems have recently been engineered to enable whole phage assembly [1], opening the door to a myriad of phage-based applications with this highly controllable and open environment. Combining this next-generation technology with the well-established field of electron microscopy (EM) allows for a unique look into the assembly dynamics of phages with a degree of precision that has been previously unachievable.

Phage-Based Therapeutics and Diagnostics

A bacteriophage is a type of virus that infects and kills bacteria. The relationship between a phage and its cognate bacterial prey constitutes the oldest predator-prey interaction on Earth, having existed for at least 1 billion years [2]. During this time, phages have evolved extreme specificity and sensitivity towards their hosts. One global problem of extreme importance, where phages can be of particular use, is tackling antimicrobial resistance (AMR). It is predicted that if the overwhelming surge of AMR is not seriously confronted over the next 30 years, then it will overtake cancer in the number of fatalities caused - with potential death tolls reaching 10 million annually. AMR is a natural process, but the misuse of antibiotics in humans and animals is accelerating the process at a dangerous rate. One approach for reducing AMR is to not disburse antibiotics in the first instance unless they are completely necessary, a second approach is to substitute the use of antibiotics with an alternative therapeutic.

Via phage deployment, both approaches can be executed by (a) deciphering whether or not antibiotics are necessary by identifying the pathogen causing the infection in a phage-based diagnostic device and (b) replacing antibiotics with Phage Therapy.

Phages offer a naturally occurring chassis that can, with the tools that synthetic biology offers, be modified and optimised for use in highly specific bacterial detection devices and therapeutic treatments. Even prior to any modifications, they offer many unique benefits over traditional approaches, including: a) High specificity and sensitivity towards their cognate host, b) Capacity to detect extremely low host presence, c) Capacity to function with impure samples under diverse or even harsh conditions, d) Discrimination between viable and incapacitated target pathogenic cells (i.e., removes false positive for diagnostics), e) Signal amplification capacity alongside signal transduction, f) Low cost, easy propagation and purification.

TXTL Phage Assembly

TXTL harnesses the endogenous transcriptional and translational machinery extracted from cells (commonly *E. coli*) and combines this cellular hardware with an energy solution and amino acid mix, allowing for the expression of genetic code in a single cell-free reaction. The ability to express and assemble phages from their isolated genome in a cell-free system allows for extensive control over when and how phages are deployed. This ground-breaking technique also gifts researchers the ability to extensively analyse phage assembly. This knowledge and increased level of control over phages is likely to supply ammunition for future phage-based therapeutics/diagnostics.

Our Main Contributions

We have extended the repertoire of phages that have been synthesised in a TXTL reaction by successfully assembling K1F phage. K1F is of particular interest because of its cognate host, *E. coli* K1 - a gram-negative pathogen, responsible for a wide range of diseases in humans, including sepsis, neonatal

*This research is supported, in parts, by the EPSRC Standard Research Studentship (DTP) EP/R513374/1, BBSRC Future Leader Fellowship (ref. BB/N011872/1) to Antonia P. Sagona, Lucidix, and EPSRC/BBSRC grant BB/M017982/1 to the Warwick Integrative Synthetic Biology Centre. We acknowledge the Midlands Regional Cryo-EM Facility, hosted at the Warwick Advanced Bioimaging Research Technology Platform, for use of the JEOL 2100Plus, supported by MRC award reference MC-PC-17136. The first author's e-mail: joseph.wheatley.1@warwick.ac.uk.

meningitis, urinary tract infections and inflammatory bowel syndrome. The novelty of our work stems from its combination of iterative TXTL phage assembly reactions and EM imaging. By taking advantage of the open and controllable nature of a TXTL system, we can gain a thorough visual insight into phage assembly and importantly, we are able to precisely quantify the start and end points of DNA expression (when the genome is added to the reaction and when the transcriptional inhibitor, rifampicin, is added to the reaction). Subsequently, these timepoints can then be accurately aligned with phage titers to produce representative data that isn't limited by variables such as TXTL reaction efficacy or phage type. This defined level of control is unattainable in vivo due to the turbulent nature of phage propagation and opaque composition of bacterial cells.

2 METHODS

TXTL K1F Phage Assembly

Cell-free reactions, using the “myTXTL” kit from Arbor Biosciences, were carried out in a volume of 12 μL for 16 hours at 29°C. Reactions were stopped at different time intervals by adding 100 $\mu\text{g}/\text{mL}$ rifampicin.

Bacteriophage Titration

The standard plaque assay was used to count K1F phage using the *E. coli* EV36 strain, an *E. coli* K12/K1 hybrid derivative with the ability to display a K1 polysaccharide capsule morphologically similar to that of *E. coli* K1 clinical isolates.

Electron Microscopy

10 μL drops of TXTL K1F assembly reaction from different time intervals were applied to the centre of the mesh and were incubated for 1 minute. The samples were removed and the mesh washed twice with 10 μL drops of water and finally negatively stained with 10 μL 2% uranyl acetate for 1 minute. Images were acquired using the Jeol 2100 transmission electron microscope.

3 RESULTS AND ANALYSIS

For the first 15 minutes after the reaction started, no observations were made (Figure 1A). Sporadic phage capsids could be seen at the half-hour mark (Figure 1B), which suggests that protein expression and capsid assembly were underway within 30 minutes of the start. Interestingly, after 45 minutes, the phage capsids had started to accumulate together and after 60 minutes this was still ongoing with increased abundance (Figure 1C and 1D). Upon analysing the TXTL K1F assembly titer data, it is revealed that the 45-minute time point aligns to a titer of 10^3 PFU/mL and the 60-minute time point aligns to a titer of 10^7 PFU/mL (Figure 1H).

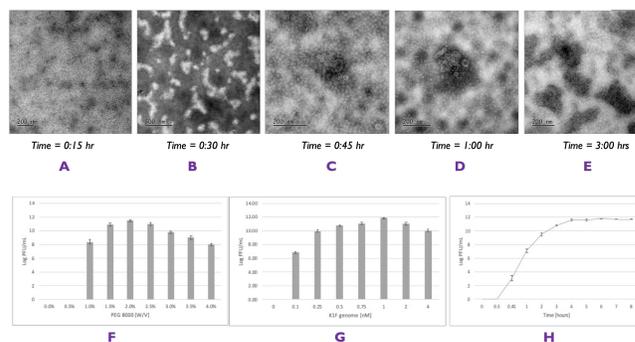


Figure 1: A-E: TXTL K1F Assembly Visualised with EM. F-G: myTXTL K1F Assembly Optimisation. H: Time Course Titer Analysis of TXTL Synthesised K1F Phage.

Our interpretation of this data would explain that the ‘capsid accumulation’ events that are displayed at 45 and 60 minutes are in fact ‘DNA packaging’ events. This would, for the first time, suggest that premature phages accumulate together in large DNA packaging events at the last stage of their development and subsequently, they become viable phages – hence the rapid shift from 0 PFU/mL at 30 minutes, to 10^3 PFU/mL at 45 minutes, to 10^7 PFU/mL at 60 minutes. After 180 minutes, the abundance of the phages can be seen to have increased, however, they appear to be much less accumulative. This suggests that the DNA packaging events are mostly complete by the 3-hour mark and that almost all phages synthesised in the reaction are viable by that time – this interpretation is supported by the fact that the amount of titer observed at the 3-hour mark is 10^{10} PFU/mL.

4 FUTURE DIRECTIONS

We plan to use this novel phage analysis platform to further validate our DNA Packaging Event theory and to investigate other aspects of K1F phage assembly as a part of our wider goal of developing phage-based therapeutic and diagnostic models. Furthermore, this platform offers a new premise for identifying and investigating genes involved in phage assembly and other purposes through iterative gene knockouts – which gives rise to a plethora of research opportunities that can be explored.

REFERENCES

- [1] RUSTAD, M., EASTLUND, A., MARSHALL, R., JARDINE, P., AND NOIREAUX, V. Synthesis of infectious bacteriophages in an *e. coli*-based cell-free expression system. *Synthetic Biology*, 126 (2017), 56144.
- [2] SHORT, F., BLOWER, T., AND SALMOND, G. A promiscuous antitoxin of bacteriophage T4 ensures successful viral replication. *Molecular Microbiology* 83, 4 (2012), 665–668.

Engineering SpyTag Bacteriophage K1F for Directional Immobilisation

Sahan B.W. Liyanagedera¹, Joseph P. Wheatley¹, Alona Yu. Biketova², Ian Hands-Portman¹, Antonia P. Sagona¹, Kevin Purdy¹, Tamas Feher², Vishwesh Kulkarni^{1*}

¹ University Warwick, ² Hungarian Academy of Sciences

{s.liyanagedera,joseph.wheatley.1,i.j.portman,k.purdy,a.sagona,v.kulkarni}@warwick.ac.uk

{alyona.biketova,fehertamas.issb}@gmail.com

1 INTRODUCTION

Pathogenic bacterial infections must be detected early and rapidly to improve the efficacy of the required therapeutic interventions. However, conventional bacterial detection methods are highly complex, require trained personnel and are time-consuming. In this context, the use of a bacteriophage, also referred to as *phage*, as a detection probe has several advantages for rapid bacterial screening including (1) extreme specificity to cognate host, (2) massive increase in progeny phage from a single infection event, and (3) simple and inexpensive large-scale production [4]. Engineering of the required genetically modified bacteriophages is becoming increasingly easier due to the development of new techniques such as CRISPR/Cas for the selection of recombinant bacteriophages [3]. We add to this rapidly evolving field by presenting the first known engineering of SpyTag K1F bacteriophage. This bacteriophage is intended for the detection of *E. coli* K1, a clinically relevant bacteria.

2 MAIN CONTRIBUTIONS

We demonstrate the first known incorporation of the SpyTag protein on to the capsid head of the K1F bacteriophage. After creating the K1F-SpyTag and K1F-GFP-SpyTag bacteriophages via homologous recombination and CRISPR/Cas mediated selection, we demonstrate their *in vitro* assembly in an *E. Coli* Cell Free *transcription/translation* (TXTL) system. We establish a method for facile directional immobilisation of SpyTag fusion proteins – and by proxy SpyTag phage – on to SpyCatcher decorated Polydimethylsiloxane (PDMS) microfluidic chips, a material commonly used in the construction of disposable *point of care* (P.O.C) diagnostic devices.

3 METHODS

Integration of SpyTag on to Minor Capsid Protein of K1F: The SpyTag and GFP-SpyTag genes were incorporated

*This research is supported, in parts, by the EPSRC Standard Research Studentship (DTP) EP/R513374/1, BBSRC Future Leader Fellowship (ref. BB/N011872/1) to Antonia P. Sagona; by EPSRC/BBSRC grant BB/M017982/1 to the Warwick Integrative Synthetic Biology Centre; by the EPSRC DTP Award, the Wellcome Trust Covid Relief Fund Grant, and the Medical and Life Sciences Research Fund grant to Sahan Liyanagedera. The first author's e-mail: s.liyanagedera@warwick.ac.uk.

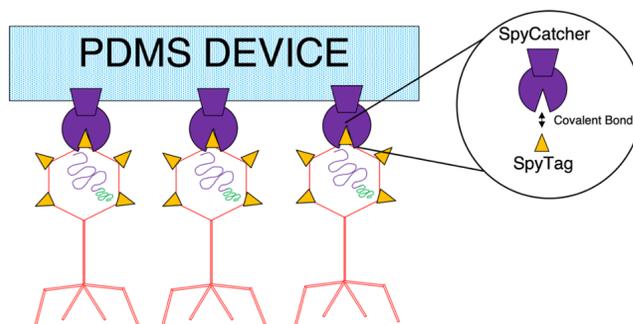


Figure 1: Concept of SpyTag-K1F Phage directionally immobilised via surface displayed SpyCatcher proteins embedded in PDMS device.

into the minor capsid protein of the bacteriophage K1F via homologous recombination [3]. To enrich the mixed population produced by homologous recombination, a selection strain containing the CRISPR/Cas machinery that targeted the wild type genome, was subject to three rounds of infection with the mixed population until stable recombinant phage were isolated via plaque assays.

PDMS Surface Decoration with SpyCatcher: Wild type Bsla and Bsla-SpyCatcher proteins were mixed in 2:1 ratio to form the spotting solution in 0.1 M carbonate buffer (pH = 9). This solution is spotted on to a hydrophobic microscope slide and allowed to dry. PDMS precursor and curing agent were mixed in a 10:1 ratio, degassed for 1 hour, and subsequently poured over the spotted microscope slide. Following curing, the PDMS chip was slowly peeled away to create PDMS chip with surface decorated SpyCatcher proteins.

Coupling SpyTag Fusions to the SpyCatcher Decorated PDMS: SpyTag fusions (GFP-SpyTag or K1F-SpyTag) were incubated on the surface of the SpyCatcher exposed PDMS chips for 10 minutes to facilitate covalent linkage between SpyTag and SpyCatcher. Coupled PDMS chips were washed 10 times in PBS and examined for retention of SpyTag-GFP via fluorescence or plaque assays for SpyTag-K1F.

Cell-Free TXTL Assembly of SpyTag-K1F: An in house *E. coli* TXTL system was set up for the assembly of bacteriophage as described in [1] but with a few minor modifications.

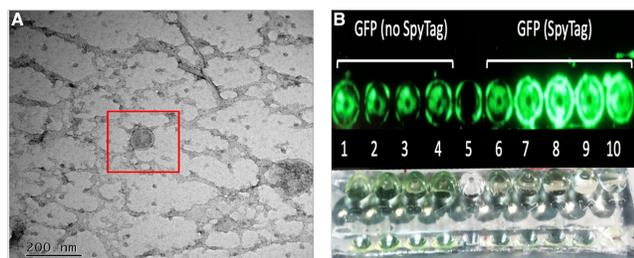


Figure 2: (A) Transmission electron microscopy image of SpyTag-K1F phage in a TXTL reaction. (B) Fluorescence image of PDMS device with immobilised SpyCatcher coupled to GFP-SpyTag. Fluorescence intensity clearly demonstrates the retention of GFP-SpyTag by surface exposed SpyCatcher proteins on the PDMS device. Non SpyTagged GFP (position 1-4) is washed away, with some fluorescence observed due to passive adsorption. Position 5 was a blank and not incubated with a coupling solution.

Quantitative characterisations of phage production in TXTL reactions were determined through plaque assays to obtain the number of plaque forming units per milliliter (PFU/mL) of TXTL reaction.

4 RESULTS AND DISCUSSION

A principle challenge in building a phage based diagnostic device stems from the difficulty in ensuring directional immobilisation of phages on the material of interest. Directional immobilisation of bacteriophages can be achieved through chemical, physio-chemical or electrostatic mechanisms between the virus and the surface in question [4]. Directional immobilisation improves stability whilst maximising the contact of bacteriophage receptor components to their target host. Nonetheless, to the best of our knowledge, no results are available for the immobilisation of bacteriophages on to PDMS, a material commonly used in the preparation of microfluidic diagnostic devices. The results in this manuscript constitute an important preliminary step in this direction.

To enable directional (tail-up) immobilisation of K1F phage, we decided to incorporate the SpyTag gene into the phage capsid head through homologous recombination and subsequent CRISPR/Cas9 selection [3]. The SpyTag protein is one half of the powerful protein conjugation pair termed the SpyTag/SpyCatcher system [5]. The system is derived by splitting the *Streptococcus pyogenes* fibro-nectin-binding protein FbaB into two functional domains, viz., SpyTag and SpyCatcher, that form a spontaneous isopeptide bond between Lysine and Asparagine residues. This ultra-strong irreversible covalent interaction is an ideal method to attach proteins to the capsid head of SpyTagged K1F phage and provides a method for its facile directional immobilisation on to surfaces containing SpyCatcher fusion proteins. We

targeted the fusion to the non essential minor capsid protein (gene10b) which is encoded when a -1 translational frame shift occurs in the 3' region of the gene encoding the major capsid protein (gene10a) [2]. In this way we were able to create two recombinant phage incorporating either SpyTag or GFP-SpyTag to the minor capsid head of K1F. Next we attempted the assembly of SpyTag-K1F in an in house Cell free TXTL system and achieved titres upwards of 10^9 PFU/mL. Cell Free assembled phage were imaged via *transmission electron microscopy* (TEM), thus demonstrating correct structure of capsids incorporating SpyTag (Figure 2 A).

5 CONCLUSION AND FUTURE DIRECTIONS

We have demonstrated a simple method to immobilise SpyTag fused proteins to SpyCatcher decorated PDMS through a method that enables control over immobilisation location on the surface of the device (Figure 2B). The physical immobilisation of SpyCatcher is achieved by spotting 1-5 μ L of high concentration of a Bsla-SpyCatcher fusion protein solution on to a microscope slide and pouring over PDMS. The use of Bsla as a fusion partner to SpyCatcher is essential to enable maximal surface exposure of SpyCatcher, as Bsla readily forms mono-layers on hydrophobic surfaces [2]. We demonstrate the ability of these SpyCatcher decorated PDMS chips to capture GFP-SpyTag as a proof of concept for the method (Figure 2B). This versatile strategy could enable the coupling of K1F-SpyTag developed in this manuscript, or any phage with SpyTag incorporated on to their capsid head, to different parts of the microfluidic device, by selective and sequential surface exposure. In doing so, a single device can be easily fabricated to detect multiple bacterial pathogens. Furthermore, we seek to build on our established methods for the assembly of SpyTag-K1F in Cell-Free TXTL, by utilising the SpyTag on the capsid head to perform affinity purification of a signal packaged phage from a TXTL reaction. Taken together these results will facilitate new generally applicable techniques for using bacteriophages as biosensors for a variety of different applications.

REFERENCES

- [1] GARENNE, D., THOMPSON, S., BRISSON, A., KHAKIMZHAN, A., AND NOIREAUX, V. The all-E. coli TXTL toolbox 3.0: new capabilities of a cell-free synthetic biology platform. *Syn. Bio.* 8, 00 (2021), 1–8.
- [2] HOBLEY, L. E. A. BslA is a self-assembling bacterial hydrophobin that coats the *Bacillus subtilis* biofilm. *PNAS* 110, 33 (2013), 13600–13605.
- [3] MOLLER-OLSEN, C., HO, S.-H., SHUKLA, R., FEHER, T., AND SAGONA, A. Engineered K1F bacteriophages kill intracellular *Escherichia coli* K1 in human epithelial cells. *Scientific Reports* 8, 17559 (2018).
- [4] ROSNER, D., AND CLARK, J. Formulations for bacteriophage therapy and the potential uses of immobilization. *Pharmaceuticals* 14, 4 (2021), 359.
- [5] ZAKERI, B., FIERER, J., CELIK, E., CHITTOCK, E., SCHWARZ-LINEK, U., MOY, V., AND HOWARTH, M. Peptide tag forming a rapid covalent bond to a protein, through engineering a bacterial adhesin. *PNAS* 109, 12 (2012), E690–E697.