# Fast and flexible simulation and parameter estimation for synthetic biology using bioscrape

**Anandh Swaminathan** [*,1], **William Poole** [†,2], **Ayush Pandey** [‡,3], **Victoria Hsiao** [4], **and Richard M Murray** [5]

**1** Ghost Locomotion, Mountain View, CA, USA **2** Altos Labs, San Francisco, CA, USA **3** Control and Dynamical Systems, California Institute of Technology, Pasadena, CA, USA **4** Amyris, Emeryville, CA, USA **5** Control and Dynamical Systems and Biology and Biological Engineering, California Institute of Technology, Pasadena, CA, USA

## Summary

In systems and synthetic biology, it is common to build chemical reaction network (CRN) models of biochemical circuits and networks. Although automation and other high-throughput techniques have led to an abundance of data enabling data-driven quantitative modeling and parameter estimation, the intense amount of simulation needed for these methods still frequently results in a computational bottleneck. Here we present bioscrape (Bio-circuit Stochastic Single-cell Reaction Analysis and Parameter Estimation) - a Python package for fast and flexible modeling and simulation of highly customizable chemical reaction networks. Specifically, bioscrape supports deterministic and stochastic simulations, which can incorporate delay, cell growth, and cell division. All functionalities - reaction models, simulation algorithms, cell growth models, partioning models, and Bayesian inference - are implemented as interfaces in an easily extensible and modular object-oriented framework. Models can be constructed via Systems Biology Markup Language (SBML) or specified programmatically via a Python API. Simulation run times obtained with the package are comparable to those obtained using C code - this is particularly advantageous for computationally expensive applications such as Bayesian inference or simulation of cell lineages. We show the package's simulation capabilities on a variety of example simulations of stochastic gene expression. We also demonstrate the package by using it to do parameter inference on a model of integrase enzyme-mediated DNA recombination dynamics with experimental data. The bioscrape package is publicly available online (Swaminathan et al., 2022) along with more detailed documentation and examples.

## Statement of need

A central theme of research in systems and synthetic biology is the quantitative predictions of the behavior of complex biological systems and efficient testing of hypotheses. Mathematical modeling and analysis tools play an integral role in such predictions and can transform the way in which we design and debug synthetic engineered biological circuits.

Cell growth and division are critical aspects of biological circuits which are typically represented as a dilution term in the model. However, in stochastic models, modeling the continuous dilution process with a stochastic and discrete degradation reaction might not be accurate. Moreover, the partitioning of molecules between daughter cells at cell division may introduce noise that is difficult to distinguish from other forms of noise (Huh & Paulsson, 2011). Therefore, modeling

---

*Co-first author
†Co-first author
‡Co-first author

39  cell growth as well as division and partitioning is important for investigating noise in gene
40  expression across a lineage of cells.

41  Regardless of simulation framework, it is necessary to first specify the values of the parameters
42  of each propensity function in the model along with the initial levels of the model species. In
43  some cases, these parameters and initial conditions are experimentally known. Often, however,
44  they have to be inferred from from biological data via a process known as parameter inference,
45  parameter estimation, or parameter identification (Sun et al., 2012). Bayesian inference
46  (Golightly & Wilkinson, 2011; Komorowski et al., 2009) is one of the most rigorous methods of
47  parameter identification. It provides a posterior distribution over the parameter space so that
48  the stochastic effects from the experimental data are modeled by the parameter distributions
49  instead of a fixed optimal point. This gives insight into the accuracy and identifiability of the
50  model. Also, such an approach allows for an easy comparison between different model classes
51  using the model evidence. The drawback of these approaches is that their implementation is
52  computationally expensive and is based on repeated forward simulations of the model within the
53  framework of Markov chain Monte Carlo (MCMC) (Golightly & Wilkinson, 2011). Therefore,
54  it is important to have the underlying simulations running as fast as possible in order to speed
55  up computation time.

56  Once a given model is fully specified, it is then important to validate the model against
57  additional biological data. In this workflow, it is often necessary to add or remove reactions
58  from the model or to perform a different type of simulation. For example, one might decide that
59  a circuit behaves too noisily for deterministic simulations and want to switch to a stochastic
60  simulation framework. If delays are playing a significant role in the dynamics, one might want
61  to incorporate previously unmodeled delays into the model.
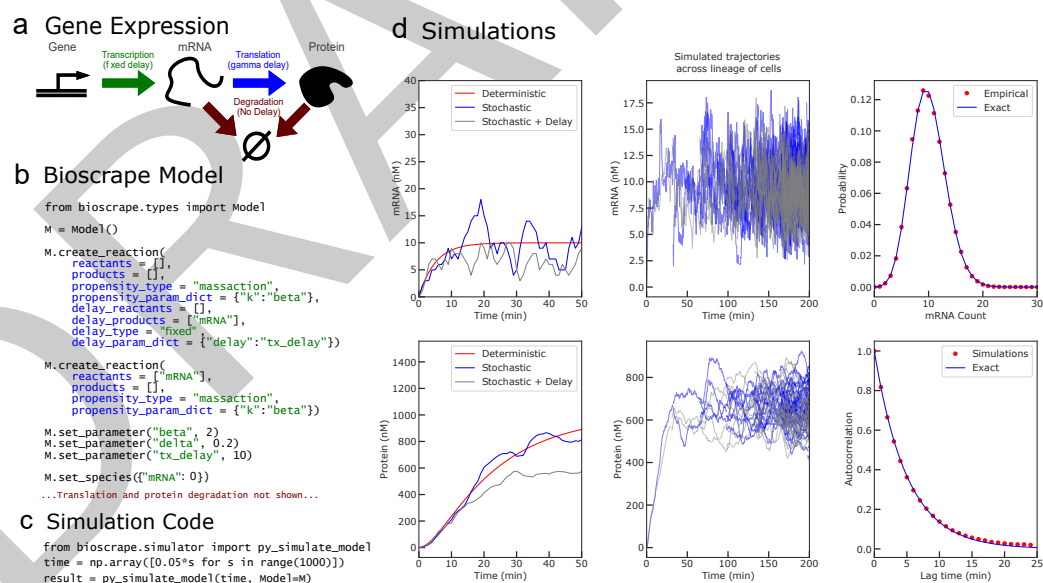


**Figure 1:** (a) A simple model of gene expression with transcription, translation, mRNA degradation, and protein degradation. The quantity of the gene encoding for mRNA is considered constant and absorbed into the transcription rate $\beta$. (b) Example Python code to construct a CRN model of gene expression using Bioscrape. (c) Models constructed via SBML or the Python API can be easily simulated with results returned as a Pandas Dataframe (McKinney, 2010). (d) Deterministic and stochastic simulations (with and without delays) using Bioscrape.The empirical probability distribution and the autocorrelation function for mRNA in the stochastic simulation matches the theoretical Poisson and exponential curve respectively

62  The result is that a very large amount of data is needed to first parameterize and then validate
63  models. The use of technologies for lab automation makes this data collection increasingly

accessible and economical. For deterministic models, this may include data collected at many different operating conditions which can be achieved with high throughput measurement techniques involving liquid handling automation (Moore et al., 2016). For stochastic models this may include large sample sizes of single cell measurements such as flow cytometry (Sachs et al., 2005; Zechner et al., 2012) and tracking single cell lineages with fluorescent microscopy (Kretzschmar & Watt, 2012). The Python API, simulation tools, and lineage module in bioscrape provide an ideal platform for such applications.

Some popular software packages that do somewhat similar tasks as the bioscrape package are MATLAB's SimBiology toolbox (MATLAB, 2016), Stochpy (Maarleveld, 2013), COPASI (Hoops et al., 2006), and Tellurium (Medley et al., 2018). The SBML simulator libRoadRunner (Somogyi et al., 2015; Welsh et al., 2023) is the state-of-the-art in deterministic and stochastic simulations of SBML models. Bioscrape simulation performance is of the same order as libRoadRunner and only around 20-30% slower while being an order of magnitude faster than other GUI-based simulators such as MATLAB and COPASI. However, the bioscrape package provides features beyond SBML simulations as it supports fully general propensity functions, provides easy-to-use parameter identification interfaces, and allows simulation of delays, and cell populations. The target audience for bioscrape includes researchers from diverse fields such as systems biology, synthetic biology, and chemical engineering. It is also aimed as an educational tool for classes on mathematical and computational biology.

## Summary of features

Figure 1 shows an example of gene expression model created and simulated stochastically and deterministically using Bioscrape.

We conclude with a list of Bioscrape features:

1. Bioscrape provides a Cython (Behnel et al., 2011) based simulator that compiles code using a C compiler to vastly increase speed. This helps assuage the computational time issues that arise in parameter estimation and stochastic simulation.
2. Kinds of possible simulations include: deterministic, stochastic, growing and dividing lineages of single cells, and stochastic simulation of delayed chemical reaction networks. A flexible easy-to-use wrapper and a Python API make it straightforward for a researcher to change their model and try simulations under diverse conditions.
3. Markov Chain Monte Carlo (MCMC) sampler based inference tools to identify parameter distributions of biological circuit models using experimental data. Bioscrape provides interfaces to easily use common biological data types such as time-series fluorescence data and flow cytometry data. The MCMC sampler is a wrapper around Python emcee (Foreman-Mackey et al., 2013).
4. Bioscrape can be used to perform local sensitivity analysis of a model to study the sensitivities of each parameter with time.

(Pandey et al., 2022) demonstrates Bioscrape's features for quantification and predictive modeling of an engineered biological system.

## Acknowledgements

111 large who have used and provided feedback on bioscrape.

## References

113 Behnel, S., Bradshaw, R., Citro, C., Dalcin, L., Seljebotn, D. S., & Smith, K. (2011).
114    Cython: The best of both worlds. *Computing in Science Engineering*, *13*(2), 31–39.
115    https://doi.org/10.1109/mcse.2010.118

116 Foreman-Mackey, D., Hogg, D. W., Lang, D., & Goodman, J. (2013). emcee: The MCMC
117    Hammer. *arXiv*, *125*, 306.

118 Golightly, A., & Wilkinson, D. J. (2011). Bayesian parameter inference for stochastic bio-
119    chemical network models using particle markov chain monte carlo. *Interface Focus*.
120    https://doi.org/10.1098/rsfs.2011.0047

121 Hoops, S., Sahle, S., Gauges, R., Lee, C., Pahle, J., Simus, N., Singhal, M., Xu, L., Mendes, P.,
122    & Kummer, U. (2006). COPASI—a complex pathway simulator. *Bioinformatics*, *22*(24),
123    3067–3074.

124 Huh, D., & Paulsson, J. (2011). Random partitioning of molecules at cell division. *Proceedings*
125    *of the National Academy of Sciences*, *108*(36), 15004–15009. https://doi.org/10.1073/
126    pnas.1013171108

127 Komorowski, M., Finkenstädt, B., Harper, C. V., & Rand, D. A. (2009). Bayesian inference of
128    biochemical kinetic parameters using the linear noise approximation. *BMC Bioinformatics*,
129    *10*(1), 343. https://doi.org/10.1186/1471-2105-10-343

130 Kretzschmar, K., & Watt, F. M. (2012). Lineage tracing. *Cell*, *148*(1-2), 33–45.

131 Maarleveld, B. G. A. B., Timo R. AND Olivier. (2013). StochPy: A comprehensive, user-
132    friendly tool for simulating stochastic biological processes. *PLOS ONE*, *8*(11), 1–10.
133    https://doi.org/10.1371/journal.pone.0079345

134 MATLAB. (2016). *Version 9.0.0 (R2016a)*. The MathWorks Inc.

135 McKinney, W. (2010). 10.25080/majora-92bf1922-00a. In S. van der Walt & J. Millman
136    (Eds.), *Proceedings of the 9th python in science conference* (pp. 51–56).

137 Medley, J. K., Choi, K., König, M., Smith, L., Gu, S., Hellerstein, J., Sealfon, S. C., & Sauro,
138    H. M. (2018). Tellurium notebooks—an environment for reproducible dynamical modeling
139    in systems biology. *PLoS Computational Biology*, *14*(6), e1006220.

140 Moore, S. J., MacDonald, J. T., Weinecke, S., Kylilis, N., Polizzi, K. M., Biedendieck, R.,
141    & Freemont, P. S. (2016). Prototype of bacillus megaterium genetic elements through
142    automated cell-free characterization and bayesian modelling. *bioRxiv*.

143 Pandey, A., Rodriguez, M. L., Poole, W., & Murray, R. M. (2022). Characterization of
144    integrase and excisionase activity in cell-free protein expression system using a modeling
145    and analysis pipeline. *bioRxiv*. https://doi.org/10.1101/2022.10.05.511053

146 Sachs, K., Perez, O., Pe'er, D., Lauffenburger, D. A., & Nolan, G. P. (2005). Causal protein-
147    signaling networks derived from multiparameter single-cell data. *Science*, *308*(5721),
148    523–529.

149 Somogyi, E. T., Bouteiller, J.-M., Glazier, J. A., König, M., Medley, J. K., Swat, M. H., &
150    Sauro, H. M. (2015). libRoadRunner: A high performance SBML simulation and analysis
151    library. *Bioinformatics*, *31*(20), 3315–3321.

152 Sun, J., Garibaldi, J. M., & Hodgman, C. (2012). Parameter estimation using metaheuristics
153    in systems biology: A comprehensive review. *IEEE/ACM Transactions on Computational*
154    *Biology and Bioinformatics (TCBB)*, *9*(1), 185–202. https://doi.org/10.1109/tcbb.2011.63

155 Swaminathan, A., Poole, W., Pandey, A., Hsiao, V., & Murray, R. M. (2022). Bioscrape:
156    Biological stochastic simulation of single cell reactions and parameter estimation. In *GitHub*
157    *repository*. GitHub. https://github.com/biocircuits/bioscrape/

158 Welsh, C., Xu, J., Smith, L., König, M., Choi, K., & Sauro, H. M. (2023). libRoadRunner 2.0:
159    A high performance SBML simulation and analysis library. *Bioinformatics*, *39*(1), btac770.

160 Zechner, C., Ruess, J., Krenn, P., Pelet, S., Peter, M., Lygeros, J., & Koeppl, H. (2012).
161    Moment-based inference predicts bimodality in transient gene expression. *Proceedings of*
162    *the National Academy of Sciences*, *109*(21), 8340–8345. https://doi.org/10.1073/pnas.
163    1200161109