

Generating Synthetic Signaling Networks for in Silico Modeling Studies

Jin Xu^a, H. Steven Wiley^b, Herbert M. Sauro^a

^a*Department of Bioengineering, University of Washington, Seattle, WA, USA*

^b*Environmental Molecular Sciences Division, Pacific Northwest National Laboratory, Richland, WA, USA*

Abstract

Predictive models of signaling pathways have proven to be difficult to develop. Reasons include the uncertainty in the number of species, the complexity in species' interactions, and the sparseness and uncertainty in experimental data. Traditional approaches to developing mechanistic models rely on collecting experimental data and fitting a single model to that data. This approach works for simple systems but has proven unreliable for complex systems such as biological signaling networks. For example, uncertainty and sparseness of the data often result in overfitted models that have little predictive value beyond recapitulating the experimental data itself. Thus, there is a need to develop new approaches to create predictive mechanistic models of complex systems. However, to determine the effectiveness of any new algorithm, a baseline model is needed to test its performance. To meet this need, we developed a method for generating artificial synthetic networks that are reasonably realistic and thus can be treated as ground truth models. These synthetic models can then be used to generate synthetic data for developing and testing algorithms designed to recover the underlying network topology and associated parameters. Here, we describe a simple approach for generating synthetic signaling networks that can be used for this purpose.

1. Introduction

In this article we describe a Julia [1] package that can be used to generate synthetic signaling networks. Such networks can be used as a basis to test novel algorithms designed to identify the topology and parameters of real signaling networks from biological data. This is particularly important for

6 building predictive models of signaling networks involved in diseases, such
7 as cancer [2, 3]. Methods for developing predictive models of biochemical
8 pathways have not changed significantly over the last 60 years. Many of the
9 commonly used techniques have been translated directly from other disci-
10 plines where systems tend to be much simpler. One approach is to collect
11 experimental data and fit it to a single model using a suitable optimization
12 algorithm [4, 5]. Given the large number of state variables and parameters
13 present in signaling network models, the availability of limited data results in
14 severe overfitting and non-identifiability of parameters [6, 7], even assuming
15 that the topology of the model is correct. Such models fail to generalize and
16 often have poor predictive value. This problem is by no means restricted to
17 just biological systems but applies to any complex system, such as weather
18 forecasting [8], climate models [9, 10], financial models [11] or hydrodynamic
19 models [12]. Given the difficulties and importance in being able to develop
20 predictive models of such systems, there is a general need for the develop-
21 ment of novel approaches that take into account the uncertainties in our
22 knowledge and ability to make measurements that will generate models with
23 the most predictive power. However, evaluating the effectiveness of any new
24 model generating or parameterization algorithm requires the availability of
25 "ground truth" models against which algorithm output can be compared.
26 One approach is to generate artificial synthetic networks that are reason-
27 ably realistic that can serve as ground truth models. Such models can be
28 used to generate artificial "experimental" data that can be used to test an
29 algorithm's ability to recover the original model and species parameters. In
30 this article, we describe a computational approach for generating synthetic
31 signaling networks.

32 Biological signaling pathways [13] are information-processing networks
33 that are used by cells to translate external signals and cues into appropri-
34 ate cell actions, such as cell growth, differentiation, movement, death or
35 metabolic activity. Signaling pathways are based on the interaction of a
36 network of proteins through a limited number of processes, such as phospho-
37 rylation, protein complex formation and targeted cleavage, degradation or
38 synthesis [14]. A cell's response to external signals is typically mediated by
39 specific receptor proteins, which carry the information into cells through a
40 series of complex steps that amplify and process the signal prior to its out-
41 put to the effector function. The amplification and signal processing steps
42 frequently use enzymatic steps, such as protein phosphorylation or proteoly-
43 sis. Signaling pathways also frequently display multiple feedback loops that

44 regulate the dynamic behavior of the network and its information processing
45 ability. In many cases the role of feedback loops and the overall topology of
46 the signaling pathways is unclear.

47 The importance of cell signaling to cancer is well known [15, 16] and
48 an understanding of how signaling networks contribute to the disease would
49 clearly be useful. If reliable predictive dynamic models of signaling pathways
50 could be developed, then more rational drug design and targeting would
51 be possible. However, developing predictive mechanistic models of signaling
52 pathways is difficult due to the large number of interactive components (in-
53 cluding potentially unknown interactions) and the sparsity of suitable data to
54 calibrate them. To address this need, we have been developing perturbation-
55 based approaches to infer the underlying topology of signaling networks [17].
56 Synthetic networks would be useful in this effort, but they must recapitulate
57 the biophysical constraints that govern signaling pathways. Thus, we set out
58 to define a general approach to generating synthetic signaling networks that
59 appropriately resemble natural ones.

60 2. Methods

61 To generate synthetic networks, we first created a list of unit processes
62 that are present in most signaling networks. These are shown in Figure
63 1 and includes catalyzed transformations (A), three binding and unbinding
64 reactions (B-D), and two phosphorylation/dephosphorylation units (E, F)
65 which include single (E) and doubly-phosphorylated (F) motifs. The three
66 binding and unbinding reactions (B-D) follow mass-action kinetic rate laws.
67 The other three units (A, E, F) follow reversible Michaelis-Menten rate laws.
68 Their corresponding rate equations are shown in Table 1. All rate constants
69 and species concentrations are assigned randomly.

70 Many proteins in signaling networks are found in complexes with other
71 proteins. As a result, the algorithm starts by defining a finite set of monomeric
72 proteins and uses these in combination with the reaction process shown in
73 Figure 1 to generate a signaling pathway that can consist of multi-protein
74 complexes. Reaction process from Table 1 are selected at random with prede-
75 fined probabilities. Each signaling pathway also has a designated input and
76 output species and the network is grown between these two points. The user
77 can also specify the number of species and the maximum number of reaction
78 units that should be included in the final network. In this way arbitrarily
79 complex networks can be generated.

Unit type	Rate laws
uni-uni (A)	$v_A = \frac{C(k_f \cdot A/K_A - k_r \cdot B/K_B)}{1 + A/K_A + B/K_B}$
uni-bi (B)	$v_B = k_f \cdot A - k_r \cdot B \cdot C$
bi-uni (C)	$v_C = k_f \cdot A \cdot B - k_r \cdot C$
bi-bi (D)	$v_D = k_f \cdot A \cdot B - k_r \cdot C \cdot D$
Single phosphorylation de-phosphorylation cycle (E)	$v_{E1} = \frac{C(k_{f1} \cdot A/K_{A1} - k_{r1} \cdot B/K_{B1})}{1 + A/K_{A1} + B/K_{B1}}$ $v_{E2} = \frac{D(k_{f2} \cdot B/K_{B2} - k_{r2} \cdot A/K_{A2})}{1 + B/K_{B2} + A/K_{A2}}$
Dual phosphorylation de-phosphorylation cycle (F)	$v_{F1} = \frac{D(k_{f1} \cdot A/K_{A1} - k_{r1} \cdot B/K_{B1})}{1 + A/K_{A1} + B/K_{B1}}$ $v_{F2} = \frac{D(k_{f2} \cdot B/K_{B2} - k_{r2} \cdot C/K_{C2})}{1 + B/K_{B2} + C/K_{C2}}$ $v_{F3} = \frac{E(k_{f3} \cdot C/K_{C3} - k_{r3} \cdot B/K_{B3})}{1 + C/K_{C3} + B/K_{B3}}$ $v_{F4} = \frac{E(k_{f4} \cdot B/K_{B4} - k_{r4} \cdot A/K_{A4})}{1 + B/K_{B4} + A/K_{A4}}$

Table 1: The rate equations of the six types of units composing the artificial random networks. The three binding and dissociation reactions (B-D) follow mass-action kinetic rate laws. The other three units (A, E, F) follow reversible Michaelis-Menten rate laws.

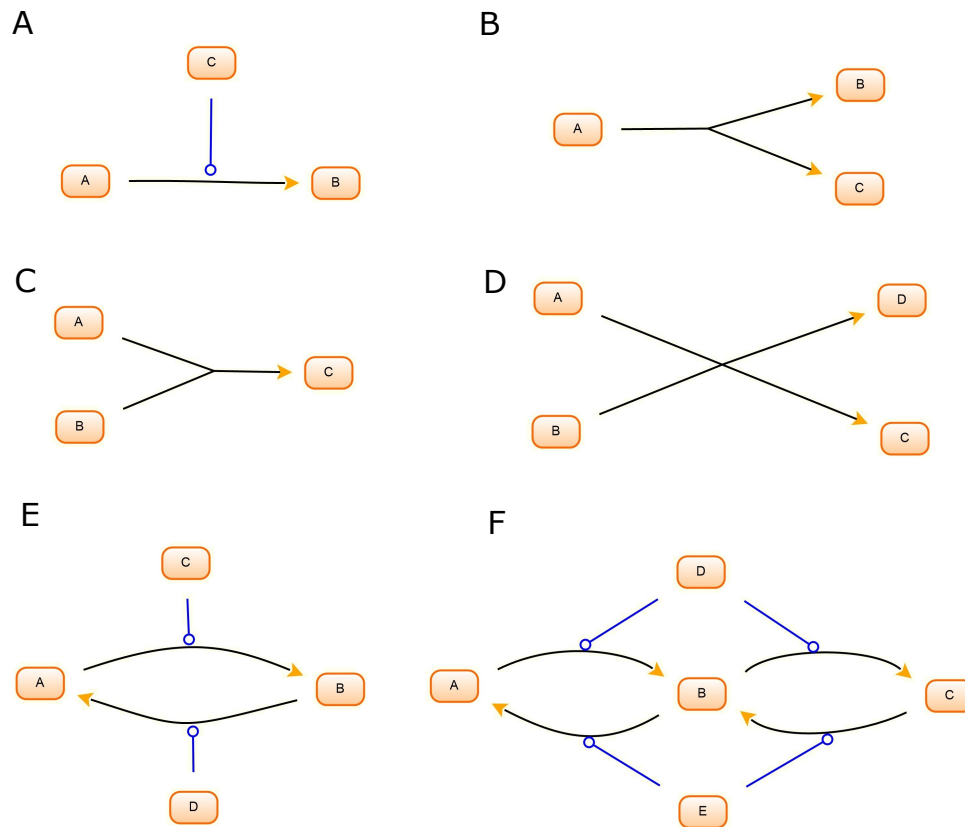


Figure 1: The synthetic random networks are composed of six types of reactions: (A) catalyzed uni-uni, (B, C) binding and unbinding reactions, (D) exchange reactions, (E) catalyzed phosphorylation-dephosphorylation, and (F) dual phosphorylation-dephosphorylation cycles.

80 However, not all initially generated networks are viable or useful. Thus,
81 the algorithm imposes structural constraints that must be followed. These
82 include requiring that all species must be connected to a path in the net-
83 work that connects the input to the output species, which ensures that no
84 isolated species exist. In addition, this prevents reaction fragments from be-
85 ing isolated from the main body. There are also dynamic constraints in the
86 network generation system. Networks that cannot reach a steady state are
87 excluded. We have also found that although some networks appear complex
88 and connected, perturbations to the input fail to propagate to the output.
89 Because such models by definition lack the ability to process information,

they were also excluded. We have also excluded networks where the input species is directly connected to the output species because those networks lack the ability to process information.

The code to generate the synthetic networks was written using the Julia Language (<https://julialang.org/>) and the simulations were done using libRoadRunner library [18] so that we could also export SBML. See appendix for details.

3. Results

Figure 2 illustrates an example of a randomly generated signaling network. It includes all six types of the chemical reaction processes in Table 1. Species are labelled S1 to S15. These include one uni-uni process: $S_{13} \rightarrow S_2$ catalyzed by S11; one bi-uni reaction: $S_{14} + S_{15} \rightarrow S_9$; one uni-bi reaction: $S_{12} \rightarrow S_8 + S_{10}$; one bi-bi reaction: $S_8 + S_{10} \rightarrow S_4 + S_{13}$; one single phosphorylation-dephosphorylation cycle: $S_7 \rightleftharpoons S_{out}$ catalyzed by S2 and S4; and a single dual phosphorylation-dephosphorylation cycle $S_3 \rightleftharpoons S_1 \rightleftharpoons S_{11}$ catalyzed by S6 and S9.

Additional random signaling networks are shown in Figure 3. As shown, all the random signaling networks have 15 species in addition to input and output species, with a limit of 15 reactions in total. The time taken to generate a single network that satisfies the constraints described in the Methods section can range from just under a minute to ten minutes for large networks (Table 2). Figure 4 illustrates some simulations where we shown how the concentrations of the output species reach a steady states. Figure 5 shows a time-course simulation of the output species, S_{out}, where there is a perturbation to the input species. All computations reported in Table 2 were done using an Intel i7 9700 processor running at 3.00 GHz with 32GB RAM using Windows 10.

Species #	15	20	25
Time (sec)	99.71 ± 127.31	247.80 ± 200.23	550.92 ± 529.40

Table 2: The time taken to generate and find a qualified random signaling network depends on the size of the network. The size of the random network is represented by the number of species involved (input and output species are not included). The errors represent the standard deviations from ten independent runs.

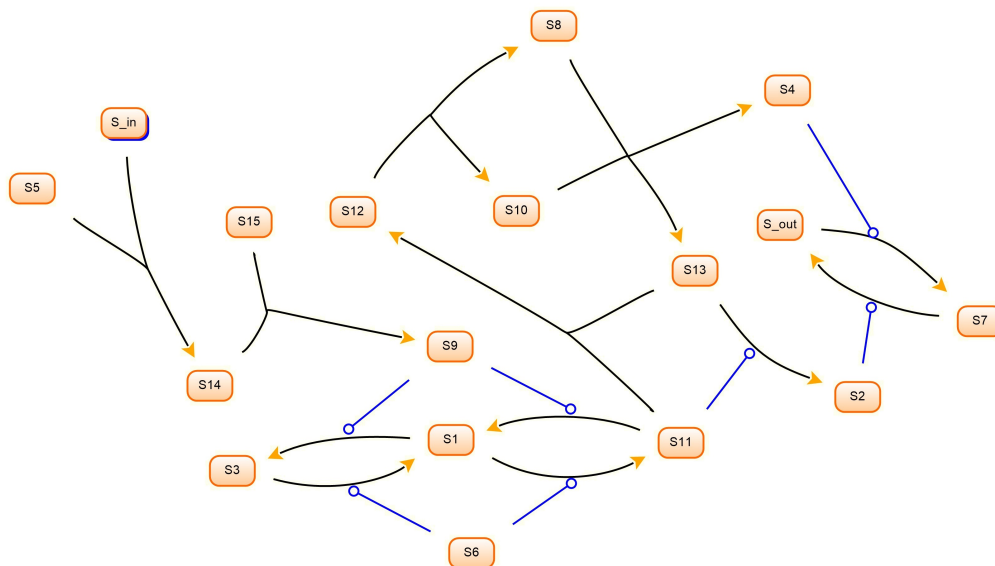


Figure 2: Random signaling networks with 15 species in addition to an input and output species (S_in and S_out) and seven catalytic interactions are shown in blue lines with a circle at the end.

117 4. Discussion

118 We have developed a Julia package that allows users to computationally
 119 generate artificial random signaling networks that can be used to test novel
 120 algorithms for parameter fitting, topology mapping and other analyses. At
 121 present we do not include any kind of sequestration in the phosphorylation
 122 cycles. Previous work has shown that sequestration of kinases and phos-
 123 phatases [19] can have a significant effect on signaling network behavior.
 124 Such effects cannot be modelled with the current version of the software.
 125 We have not examined whether our generated networks have similar graph
 126 metrics to networks found in nature, which is another area that should be in-
 127 vestigated. A more thorough investigation into the potential dynamics of our
 128 random networks needs to be carried out in order to investigate how the be-
 129 havior compares to natural networks. We have also not imposed biophysical
 130 constraints on species parameters. Because network behavior is a function of
 131 both topology and species parameters [20], such constraints could reduce the
 132 number of topological configurations that could produce viable networks. Fi-
 133 nally, further analysis of which constraints must be either imposed or relaxed

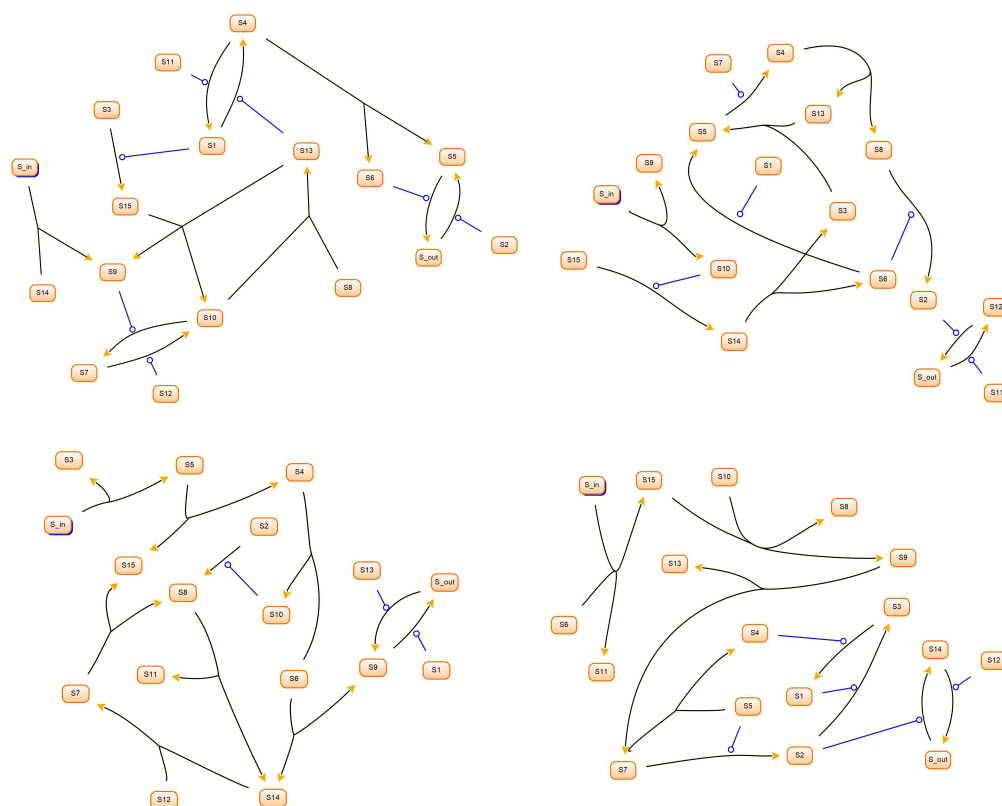


Figure 3: Four additional examples of random signaling networks with 15 species in addition to input and output species, S_{in} and S_{out} .

134 to generate realistic network behavior could reveal important principles that
135 real biological networks follow during their evolution.

136 5. Author contributions

137 JX wrote the code, generated the results and wrote the initial article
138 draft. HSW assisted in writing the article. HMS conceived the idea and
139 assisted in writing the article.

140 6. Acknowledgement

141 This work was financially supported by the National Institute of Gen-
142 eral Medical Sciences and National Cancer Institute under grant number

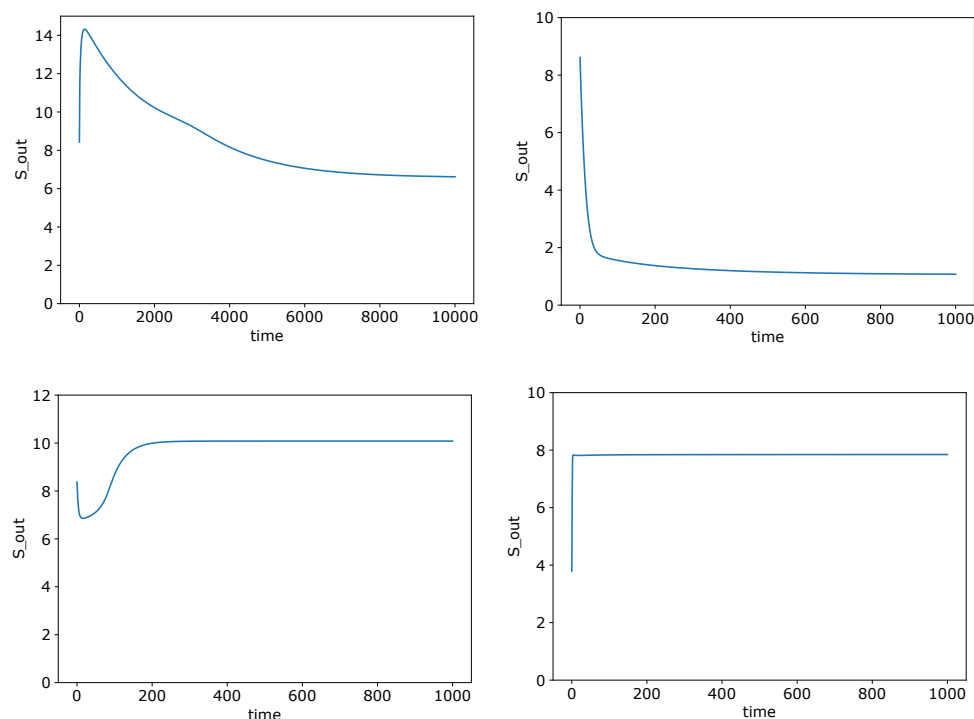


Figure 4: Time-course simulations for four randomly generated networks made up of 15 species in addition to input and output species. The plots show how the concentrations of the output species (S_{out}) reach their steady states.

143 GM12303201 and U01CA242992. The authors are responsible for the con-
 144 tent, which does not necessarily represent the opinions of the National In-
 145 stitutes of Health. We thank Luke Y. Zhu for assisting with development of
 146 some of the early code.

147 Appendix A. Code availability

148 The code to generate the artificial signaling random networks is avail-
 149 able at [https://github.com/sys-bio/aritificial_random_signaling_](https://github.com/sys-bio/aritificial_random_signaling_network)
 150 **network**. This package was implemented in Julia 1.2 on Windows 10. To use
 151 the package, first install Julia by going to the website: [https://julialang.](https://julialang.org/)
 152 **org/**,

153 Once you have the Julia console open, make sure you have **StatsBase**
 154 installed by typing `using Pkg` followed by `Pkg.add("StatsBase")`. Note

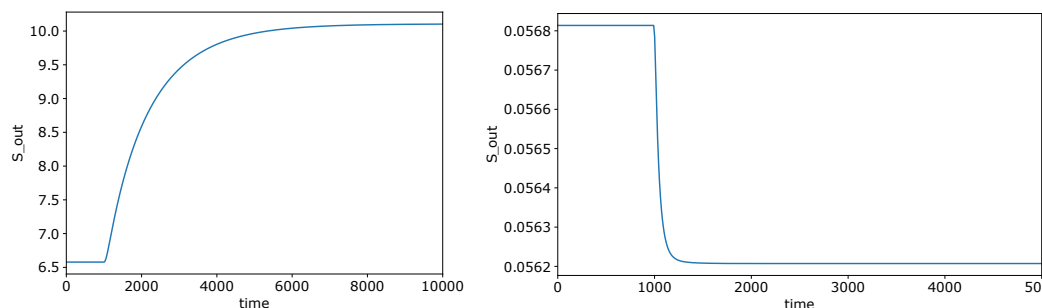


Figure 5: Time-course simulations of concentrations for the output species, S_{out} , as a result of a perturbation to the input species. The two subplots correspond to different artificial random networks with 15 species in addition to input and output species.

155 there shouldn't be a space between the `add` and the first parenthesis. This is
156 a package required for random number generation.

157 Next, download all the files from Github, [https://github.com/sys-bio/](https://github.com/sys-bio/artificial_random_signaling_network)
158 `artificial_random_signaling_network`, into one folder. At the Julia
159 console type:

160 `include("pathto\\Ground_truth_generation.jl")` to run the main script.

161 Note that `pathto` is the path you where you saved the network generation
162 scripts to, and the double backslash is to avoid the Julia misinterpreting
163 the backslash as a control character. A random signaling network will be
164 generated in SBML [21] format called `sampleNetwork.xml`. There are some
165 configuration settings in the Julia script file `Ground_truth_generation.jl`.
166 These include:

- 167 1. The number of species, `nSpecies`, and maximum number of reactions
168 `nRxns_limitation`.
- 169 2. Randomly assigned ranges for species concentrations and rate constants
170 can be modified via `rnd_species` and `rnd_parameter`.
- 171 3. `concentration_perturb` can be used to set the factor that perturbs
172 the concentration at the input species. Default is set to 2.
- 173 4. The number of random networks to generate can be set by changing
174 the variable `sampleSize`.

175 The file `roadrunner_c_api.dll` is the dynamic library for libRoadRunner
176 that is used to provide SBML simulation support [17].

References

- [1] Bezanson, Jeff & Edelman, Alan & Karpinski, Stefan & Shah, Viral. (2015). Julia: A Fresh Approach to Numerical Computing. 10.1137/141000671.
- [2] Logue, Jeremy & Morrison, Deborah. (2012). Complexity in the signaling network: Insights from the use of targeted inhibitors in cancer therapy. *Genes & development*. 26. 641-50. 10.1101/gad.186965.112.
- [3] Kuperstein, Inna & Bonnet, E & Nguyen, Hien-Anh & Cohen, David & Viara, Eric & Grieco, Luca & Fourquet, S & Calzone, Laurence & Russo, Christophe & Kondratova, Maria & Dutreix, Marie & Barillot, Emmanuel & Zinovyev, Andrei. (2015). Atlas of Cancer Signalling Network: A systems biology resource for integrative analysis of cancer data with Google Maps. *Oncogenesis*. 4. 10.1038/oncsis.2015.19.
- [4] Jaqaman, Khuloud & Danuser, Gaudenz. (2006). Linking data to models: Data regression. *Nature reviews. Molecular cell biology*. 7. 813-9. 10.1038/nrm2030.
- [5] Lillacci, Gabriele & Khammash, Mustafa. (2010). Parameter Estimation and Model Selection in Computational Biology. *PLoS computational biology*. 6. e1000696. 10.1371/journal.pcbi.1000696.
- [6] Hengl, Stefan & Kreutz, Clemens & Timmer, J & Maiwald, Thomas. (2007). Data-Based Identifiability Analysis of Nonlinear Dynamical Models. *Bioinformatics (Oxford, England)*. 23. 2612-8. 10.1093/bioinformatics/btm382.
- [7] Oana, Chis & Banga, Julio & Balsa-Canto, Eva. (2011). Structural Identifiability of Systems Biology Models: A Critical Comparison of Methods. *PloS one*. 6. e27755. 10.1371/journal.pone.0027755.
- [8] Bauer, Peter & Thorpe, Alan & Brunet, Gilbert. (2015). The quiet revolution of numerical weather prediction. *Nature*. 525. 47-55. 10.1038/nature14956.
- [9] Pierce, David & Barnett, Tim & Santer, Benjamin & Gleckler, Peter. (2009). Selecting Global Climate Models for Regional Climate Change

- 208 Studies. Proceedings of the National Academy of Sciences of the United
209 States of America. 106. 8441-6. 10.1073/pnas.0900094106.
- 210 [10] Flato, Gregory & Marotzke, J. & Abiodun, Babatunde & Braconnot,
211 Pascale & Chou, Sin Chan & Collins, W. & Cox, Peter & Driouech,
212 Fatima & Emori, S. & Eyring, V. & Forest, Chris & Gleckler, P. &
213 Guilyardi, Eric & Jakob, C. & Kattsov, V. & Reason, Chris & Rum-
214 mukainen, M.. (2013). Evaluation of climate models.
- 215 [11] Seneta, Eugene. (2004). Fitting the variance-gamma model to
216 financial data. Journal of Applied Probability. 41. 177-187.
217 10.1017/S0021900200112288.
- 218 [12] Alexander, R. & Clarke, C. & Pringle, J.. (2006). Photoevapora-
219 tion of protoplanetary discs – I. Hydrodynamic models. MNRAS. 369.
220 10.1111/j.1365-2966.2006.10293.x.
- 221 [13] Basson, M.. (2012). Signaling in Cell Differentiation and Morphogen-
222 esis. Cold Spring Harbor perspectives in biology. 4. 10.1101/cshper-
223 spect.a008151.
- 224 [14] Klipp, Edda & Liebermeister, Wolfram. (2006). Klipp, E. & Lieber-
225 meister, W. Mathematical modeling of intracellular signaling pathways.
226 BMC Neurosci. 7, S10. BMC neuroscience. 7 Suppl 1. S10. 10.1186/1471-
227 2202-7-S1-S10.
- 228 [15] Somogyi, Endre & Bouteiller, Jean-Marie & Glazier, James & König,
229 Matthias & Medley, Kyle & Swat, Maciej & Sauro, Herbert. (2015).
230 libRoadRunner: A High Performance SBML Simulation and Analysis
231 Library. Bioinformatics. 31. 10.1093/bioinformatics/btv363.
- 232 [16] Sever, Richard & Brugge, Joan. (2015). Signal Transduction in Can-
233 cer. Cold Spring Harbor Perspectives in Medicine. 5. a006098-a006098.
234 10.1101/cshperspect.a006098.
- 235 [17] Choi, Kiri & Hellerstein, Joesph & Wiley, H Steven & Sauro Herbert M.
236 (2018). Inferring Reaction Networks using Perturbation Data. BioRxiv.
237 1-9. 10.1101/351767.
- 238 [18] Somogyi, Endre & Bouteiller, Jean-Marie & Glazier, James & König,
239 Matthias & Medley, Kyle & Swat, Maciej & Sauro, Herbert. (2015).

- libRoadRunner: A High Performance SBML Simulation and Analysis
Library. *Bioinformatics*. 31. 10.1093/bioinformatics/btv363.
- [19] Markevich, Nikolai & Hoek, Jan & Kholodenko, Boris. (2004). Sig-
naling switches and bistability arising from multisite phosphorylation
in protein kinase cascades. *The Journal of cell biology*. 164. 353-9.
10.1083/jcb.200308060.
- [20] Klinke, David J. (2010). Signal Transduction Networks in Cancer: Quan-
titative Parameters Influence Network Topology. *Cancer Research*. 70.
1773-82. 10.1158/0008-5472.CAN-09-3234.
- [21] Hucka, Michael & Finney, A. & Sauro, Herbert & Bolouri, H. & Doyle,
John & Kitano, Hiroaki & Arkin, Adam & Bornstein, B.J. & Bray,
D. & Cornish-Bowden, Athel & Cuellar, A.A. & Dronov, S. & Gilles,
E.D. & Ginkel, M. & Gor, V. & Goryanin, Igor & Hedley, W.J. &
Hodgman, Charlie & Hofmeyr, Jan-Hendrik & Wang, J.. (2003). The
systems biology markup language (SBML): a medium for representation
and exchange of biochemical network models. *Bioinformatics*. 19. 524-
531. 10.1093/bioinformatics/btg015.