

# 1 Sensitivity analysis of Repast 2 Computational Ecology models with 3 R/Repast

4 Antonio Prestes García<sup>1</sup> and Alfonso Rodríguez-Patón<sup>1</sup>

5 <sup>1</sup>Departamento de Inteligencia Artificial, Universidad Politécnica de Madrid, Campus de  
6 Montegancedo s/n, Boadilla del Monte, Madrid, Spain

## 7 ABSTRACT

8 Computational ecology is an emerging interdisciplinary discipline founded mainly on modeling and  
9 simulation methods for studying ecological systems. Among the existing modeling formalisms, the  
10 individual-based modeling is particularly well suited for capturing the complex temporal and spatial  
11 dynamics as well as the nonlinearities arising in ecosystems, communities or populations due to individual  
12 variability. In addition, being a bottom up approach, it is useful for providing new insights on the local  
13 mechanisms which are generating some observed global dynamics. Of course no conclusions about  
14 model results could be taken seriously if they are based on a single model execution and they are  
15 not analyzed carefully. Therefore, a sound methodology should always be used for underpinning the  
16 interpretation of model results. The sensitivity analysis is a methodology for quantitatively assessing  
17 the effect of input uncertainty in the simulation output which should be incorporated compulsorily to  
18 every work based on in silico experimental setup. In this paper we present R/Repast a GNU R package  
19 for running and analyzing Repast Symphony models accompanied by two worked examples on how to  
20 perform global sensitivity analysis and how to interpret the results.

21 Keywords: Individual-Based Modeling, Sensitivity analysis, Repast, Computational Ecology, Systems  
22 Biology

## 23 INTRODUCTION

24 The computational ecology is a relatively young field which relies extensively on mathematical com-  
25 putational methods and models for studying ecological and evolutionary processes. It is based on the  
26 construction of predictive and explanatory models as well as the quantitative description and analysis  
27 of ecological data Helly et al. (1995) Petrovskii et al. (2012). The continuous growth of computational  
28 power available for and end users, the existence of tools and the constant increment of empirical data  
29 available, makes viable for many scientists to develop and simulate tremendously complex models from  
30 their desktops. In addition, the intrinsic characteristics of ecological processes, maxim their temporal and  
31 spatial scale Dieckmann et al. (2000), converts the task of carrying out controlled experiments a physical  
32 impossibility. Hence, in most cases the only feasible alternative is to simulate the process in order to make  
33 experiments spanning the full length of ecological and evolutionary scales. The computational ecology  
34 has its roots from the successful results achieved from mathematical ecology which has proven to be an  
35 essential tool for understanding the complexities which arise from ecological interactions.

36 It is a widely accepted that simple models with a small number of state variables and parameters  
37 provide best generalizations than the complex ones Smith (1974) Evans et al. (2013) with a clear distinction  
38 between simulation models and theories as separate entities handling different kind of problems. It has  
39 been recently questioned the correctness of the idea the simple models lead to generality in ecology Evans  
40 et al. (2013). We believe that the parsimony principle must always be taken into account when developing  
41 models, but this has a different meaning depending on the modeling formalism we are using. Simplicity  
42 does not have the same meaning when the referred modeling formalism is a deterministic ODE or when it  
43 is applied to Agent-based modeling, as long as every modeling techniques has its own idiosyncrasy and  
44 constraints. The Agent-based modeling is a flexible and versatile abstraction where the whole system  
45 under study is described or formalized by its component units, which facilitates a more natural description

46 of a system and the comprehension of individual properties leading to the emergent phenomena Bonabeau  
47 (2002).

48 The Agent-based models (AbM) are much more fine-grained than their whole-population aggregated  
49 counterpart and as consequence they tend to be more complex requiring more equations, parameters and  
50 processes in order to represent the same phenomenon. That is not intrinsically a problem or a quality but  
51 simply a constraint imposed by the modeling formalism in use and it is up to the modelers to find the  
52 correct tradeoff between the purpose of the model and the level of details which should be part of the  
53 model structure.

54 The AbM is being established progressively as a main-stream and valuable tool for modeling complex  
55 adaptive systems in many distinct areas of knowledge, ranging from social science, economics to any  
56 flavor of computational and systems science such as biology, ecology and so on Grimm and Railsback  
57 (2005). The reason is, amongst other things, the relative ease with which detailed structural information  
58 can be incorporated into a model without the constraints of other methodologies Hellweger and Bucci  
59 (2009). Nonetheless, the possibility of incorporating many details comes with the cost of models with a  
60 high complexity level, containing many rules and parameters for which the exact values are, in many cases,  
61 hard or impossible to determine experimentally, that is what is known as parameter uncertainty. When  
62 used in the context of ecological systems the Agent-based modeling is also known as Individual-based  
63 modeling (IbM) Grimm and Railsback (2005).

64 The distinctive aspect defining what is an IbM is that individuals are represented by discrete entities  
65 and they also have a property or state variable which are unique in the population being simulated Berc  
66 (2002). Hence IbM is a valuable abstraction for simulating populations, communities or ecosystems  
67 capturing the individual variability, randomness and their complex dynamics. It is a bottom-up approach  
68 where the system under study is modeled using mechanistic explanations on the interacting system parts  
69 Ferrer et al. (2008). Therefore, the global behavior shown by the system as a whole, is an emergent  
70 property derived from the local rules defining the individuals. That is particularly useful testing different  
71 hypothesis or phenomenological explanations for the individual processes in order to verify which of  
72 them are producing the global observed behavior Pascual (2005). Moreover, differently from aggregate  
73 models, it is customary that IBM have a large number of state variables and parameters which in most  
74 cases are hard or directly impossible to elucidate experimentally leading to many levels of uncertainty in  
75 this kind of models. In order to tackle with the uncertainty and for making robust predictions, we have  
76 to use a sound methodology for applying what-if analysis to check how stable are the model outputs  
77 when varying the input parameters Thiele et al. (2014). There exist a large set of mathematical tools for  
78 analyzing the model output which are known generically as sensitivity analysis. Normally applying these  
79 techniques are cumbersome, requiring a lot of effort from modelers, hindering the throughout analysis of  
80 computational models.

81 According to Thiele et al. (2014) most of Individual-based models published tends to omit the  
82 systematic analysis of model output, mainly because modelers normally do not have the specific knowledge  
83 to implement the required methods. Therefore, it seems to be clear, that the availability of simple and  
84 user friendly tools for experiment design and analysis would greatly help modelers to improve the formal  
85 quality of their models.

86 In other scientific fields, which are strongly rooted on an extensive experimentalism, is practically  
87 impossible to conduct any kind of research without a well-designed experimental setup and a further  
88 statistical analysis and hypothesis test. Perhaps the reasons are that these experimental fields already have  
89 a complete and mature toolbox for design and evaluation of experiments Little and Hills (1978) Myers  
90 and Well (1995) leaving no room for deviation from these standards. On the other hand, in silico based  
91 experiments are still on early stage and verification and validation procedures are not well established yet.  
92 In addition, the real value of a computational model depends much on the ability of other researchers  
93 to reproduce and enhance the results elsewhere; in other words, results must be reproducible. Hence, in  
94 order to achieve reproducibility, research methods should be stated clearly and should preferentially being  
95 backed by standard methods and software tools.

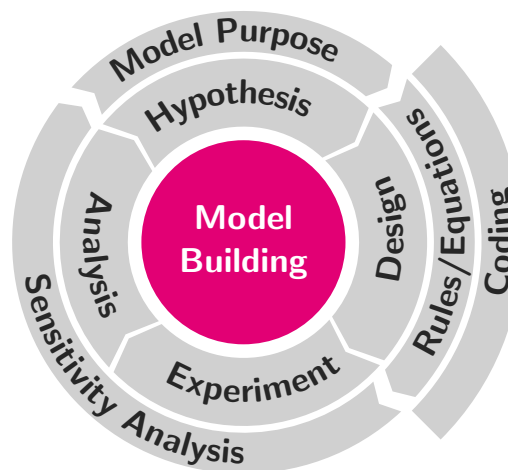
96 Bearing this in mind we introduce R/Repast a GNU R package for running Repast North et al. (2013b)  
97 models from GNU R environment as well as for carrying out global sensitivity analysis on the model  
98 results. In the following sections we will contextualize the problem providing a basic background for  
99 understanding what is being addressed in this work and we will also provide a basic description about  
100 the package functionalities. Finally, we will show three worked examples on how the package can help

101 modelers to make the conclusions drawn from model results much more robust. The first example explores  
102 the basic aspects of bacterial conjugation process. The second is an individual-based implementation of  
103 the classic predator-prey model enclosed as part of the standard Repast Symphony distribution. Finally,  
104 the last example was developed ex professo for this work and it is an instance of common pool problem in  
105 the context of two plasmids "sharing" the genes required for the expression of conjugative system.

## 106 BACKGROUND

### 107 Model development

108 Model development is an iterative and objective driven activity and the first step required to develop a  
109 model is having a clear and ideally unambiguous statement about the model purpose. Therefore, every  
110 experimental study carried out using modeling and simulation should follow the experimental life cycle  
111 based on the successive sequence of four cyclic steps, starting from (a) Conjecture, which defines the  
112 model purpose and why the model is being developed; (b) Design phase where the model is translated to  
113 some runnable implementation; (c) Experiment step which means the execution of model following a  
114 well-established plan oriented to confirm or reject the initial conjecture and finally the (d) Analysis step  
115 where the data generated in the previous step is analyzed with a sound methodology which, hopefully will  
116 generate new insights, uncover model flaws and iteratively improve the initial conjecture and design Box  
117 and Draper (1987). A simple graphical representation of these four iterative steps is shown in Figure 1.



**Figure 1.** The iterative model development life cycle. This figure shows the relationship between the modeling phases and their associated tasks when applied to an individual-based model.

118 Part of design phase consist in convert the model equations and rules to a computer code implementa-  
119 tion. Currently there are several frameworks available for developing individual-based models. These  
120 frameworks are designed to address some specific requirement such as usability Tisue and Wilensky  
121 (2004), flexibility or scalability North et al. (2013b); Luke et al. (2005) or support to multiple paradigm,  
122 such as AnyLogic Emrich et al. (2007). Certainly, the most widespread framework in ecological modeling  
123 is NetLogo Tisue and Wilensky (2004) which is considered to provide an easier development environment  
124 based on extensions to Logo paradigm especially suited for those which are not much familiar with  
125 modern programming languages. One of the main drawbacks of NetLogo is the scalability. NetLogo tends  
126 to show some performance issues when simulating a large number of agents. On the other hand, Repast  
127 Symphony framework has a steep learning curve but provides a fast and flexible java-based environment  
128 with many interesting features for simulating large scale computational ecology models. These features  
129 include, amongst others things, the integration with Weka, exporting the model output to R environment,  
130 support for running distributed batch simulations and some built-in facilities for parameter sweeping  
131 North et al. (2013b). Finally, Mason is, in some extent, very similar to Repast but less mature than it

132 is; It has been designed focusing on providing faster execution speeds Luke et al. (2005). The only of  
 133 these frameworks providing integrated sensitivity analysis capabilities is AnyLogic, the other frameworks:  
 134 NetLogo, Repast and Mason which are all free software do not have built in support to sensitivity analysis.

135 The Repast framework is widely used in many different fields for building individual-based simulation  
 136 models of dynamic processes Watkins et al. (2015) Gutfraind et al. (2015) Tack et al. (2015). In addition,  
 137 Repast also has a framework for high performance computing using the C++ programming language with  
 138 similar conceptual entities as those found in Repast-java. Repast also has support for running GNU R code  
 139 R Core Team (2015) Crawley (2007) from inside the user interface but until now it has not been feasible  
 140 to run Repast models from R environment for controlling model in order to implement experimental  
 141 designs, calibration, parameter estimation and sensitivity analysis, therefore hindering a throughout and  
 142 comprehensive validation of Individual-based models developed using Repast Simphony.

### 143 Sensitivity analysis

144 Because of sensitivity analysis is a broad and complex subject, a throughout discussion would be lengthy  
 145 and out of the scope of this work. Instead we will try to provide a more amenable and practical approach  
 146 keeping the discussion at a general level but rigorous enough to let the practitioners gain the knowledge  
 147 required to understand, apply and interpret the results. For a more detailed review please refer to Saltelli  
 148 et al. (2004) Pianosi et al. (2016). It is interesting to start the discussion providing the exact meaning  
 149 of some the many expressions which are used commonly in the analysis of models. There are several  
 150 terms used in the context of sensitivity analysis for which is important to provide the formal meaning. For  
 151 instance, the jargon of sensitivity analysis includes model calibration and parameter estimation which  
 152 many times are used as they were equivalent, even though they are different objectives. Other terms such  
 153 as uncertainty analysis, omitted variable bias, objective function or cost function are also important part  
 154 of SA lexicon.

155 Generally speaking, the objective of SA is to understand the effect of varying input factors on the  
 156 model output Saltelli et al. (2004). Under this very general statement we have a wide range of methods  
 157 and techniques which are suitable for distinct kinds of models. In order to improve this definition, it is  
 158 convenient to provide a more formal definition to the entity which is the target of SA: the model. Formally  
 159 speaking, a model is a functional relation between a number  $k$  of input factor, also called independent  
 160 or predictor variable and the output variable, sometimes referred as dependent or response variable Box  
 161 and Draper (1987) as depicted by the expression  $\eta = f(x_1, x_2, \dots, x_k)$ , being  $\eta$  is the average value of  
 162 response variable considering any specific setting for the input factors  $x_i$ . Therefore the value of a single  
 163 model run is given by  $y = f(x_1, x_2, \dots, x_k) + \varepsilon$ , where  $\varepsilon$  is difference between the value of  $y$  and the  
 164 expected value  $E(y) = \eta$ . The error  $\varepsilon$  is consequence of stochasticity introduced by design in the structure  
 165 of model to capture the population variability. Finally, recognizing that most real world models usually  
 166 have more than one response variable, the structure of an individual-based model  $M$  can be generalized  
 167 for  $n$  outputs as can be seen below

$$M = \begin{cases} y_1 = f_1(x_1, x_2, \dots, x_k) + \varepsilon \\ y_2 = f_2(x_1, x_2, \dots, x_k) + \varepsilon \\ \vdots \\ y_n = f_n(x_1, x_2, \dots, x_k) + \varepsilon \end{cases}$$

168 Therefore, being  $y_i$  some output of model  $M$ , the **model calibration** process consists in comparing  
 169 these outputs to some reference values Zeigler et al. (2000) which are normally, in the case of ecological  
 170 or biological studies, experimental or observed data. The target of calibration process is minimizing  
 171 the discrepancies between simulated and reference values. The function used for computing how far  $y_i$   
 172 output is from the reference values is known as **objective function** or **cost function**. There are many  
 173 options for implementing the objective function and the only requirement is that the return of objective  
 174 function should be inversely proportional to the quality of fit, being zero the return value for the perfect  
 175 fit. Common implementations for objective function are based on the definition of acceptable ranges,  
 176 least squares or even a combination of both. For instance, let  $y_i$  be the output of some hypothetical model  
 177  $M$ , assuming this variable represents the net reproductive rate  $R_0$ . The reference values  $R_v$  for the output  
 178 variable must fall between 0.8 and 1.2, hence any  $y_i$  value within this interval is considered to have a  
 179 perfect fit, bearing this in mind the cost function could be given by the following expression

$$C(y_i) = \begin{cases} 0, & \text{if } 0.8 \leq y_i \leq 1.2 \\ 1, & \text{otherwise} \end{cases}$$

180 That is what is known as categorical calibration criteria Thiele et al. (2014). The main drawback of this  
 181 approach is that it does not provide any information about how far is the response value from the reference  
 182 value. A better alternative is to apply some distance function  $d(y_i, R_y)$  to the output and the reference  
 183 values, even standalone or in combination with categorical calibration. The most commonly used metric is  
 184 some of the multiple forms of squared deviation but any distance function can be alternatively employed  
 185 as long as two properties hold:  $d(y_i, R_y) = 0$  if  $x_i$  and  $R_y$  are equal and  $d(y_i, R_y) > 0$  when  $x_i$  and  $R_y$  are  
 186 not equal.

Whilst calibration is a general term, meaning fundamentally the comparison of some value to a  
 reference value, the term **parameter estimation** has a more subtle and specific goal. The parameter  
 estimation is normally considered an inverse problem because the objective is finding the values for the  
 model parameters providing the best adjustment to the reference values. In other words, knowing the  
 expected values for response variable the target is estimating the suitable values for the model parameters.  
 Usually the terminology parameter refers to the constants which are part of models with clear distinction  
 between parameters and independent variables, Beck and Arnold (1977), for instance in the growth  
 differential equation shown below

$$dN/dt = rN$$

187 the model parameter would be only the growth rate  $r$  and the independent variable the time, but for the  
 188 purpose of this work we consider indistinctly the model constants and independent variables as being  
 189 parameters.

190 The two main objectives of sensitivity analysis are understanding how robust are the model results  
 191 considering the existing uncertainties and quantifying the effect of input factors on the variance of output  
 192 Saltelli et al. (2004) Pianosi et al. (2016) Law (2005). The intrinsic characteristics of individual-based  
 193 models which relies on mechanistic descriptions favors the production of models with many sub-processes,  
 194 state variable and parameters. The design is normally based on incomplete knowledge resulting in several  
 195 levels of uncertainties in the model parameters, in the model response variables and in the model structure  
 196 itself. The model structure is also related to the identifiability problem where not all model parameters can  
 197 be uniquely estimated. The sensitivity analysis can be also used for assess the effect of model structure on  
 198 the output considering the alternative model implementations as being another parameter. This can be  
 199 useful for analyzing the **omitted variable bias**, which basically means that some parameter of model can  
 200 be over or under-estimated because another important parameter was not included in the model structure.  
 201 The sensitivity analysis can be carried out letting the parameters varying over the full range of parameter  
 202 space or restricted to a small region close to the average value, respectively referred as global sensitivity  
 203 analysis and local sensitivity analysis. Sensitivity analysis can also be performed varying one factor at  
 204 time (OAT) leaving all others fixed or varying all factors at the same time (AAT). The application of  
 205 second method is required in order to capture interaction between parameters and non-linear effects.

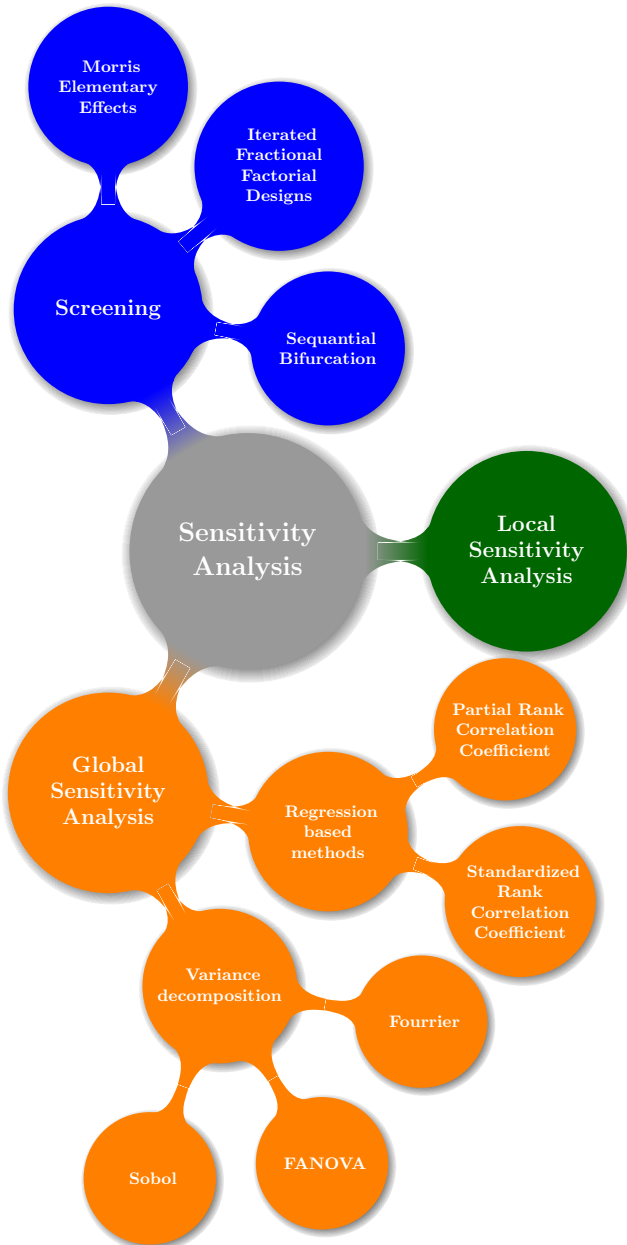
The central point of SA methodology is the estimation of sensitivity indices or coefficients. The  
 sensitivity coefficients allow the quantitative comparison of the contributions from distinct parameters to  
 the model output. In its classical form Beck and Arnold (1977) the sensitivity indices are defined as the  
 first derivative with respect to some model parameter  $x_i$ . Considering the general model  $y = f(X)$ , being  
 $X$  the parameter vector of size  $k$ , the sensitivity index  $S_i$  is given by

$$S_i = \frac{\partial Y}{\partial x_i}$$

206 It is also important to take into account that the partial derivatives can have different units, hence can be  
 207 necessary to scale them in order to make them comparable. In this approach, input factors are perturbed  
 208 one-at-time, being that measure of sensitivity suitable for local SA Pianosi et al. (2016).

209 Several methods for estimate sensitivity indices which are adequate for global sensitivity analysis are  
 210 available, such as meta-modeling approach Happe et al. (2006), correlation based methods, regression  
 211 based methods, Fourier Amplitude Sensitivity Test (FAST) Xu and Gertner (2011), for a more in depth  
 212 discussion please refer to Thiele et al. (2014) Saltelli et al. (2004) Saltelli (2008) Pianosi et al. (2016) Pujol  
 213 et al. (2015). The Figure 2 show how are related the different methods for assessing the importance of input

214 factors in simulation models, also including screening techniques Bettonvil and Kleijnen (1996) Andres  
 215 and Hajas (1993). In this work we will focus on those methods based on the variance decomposition  
 216 which are suitable for a wide range of situations, including those which are commonly found in individual-  
 217 based models, such as non-linear mappings between input factors and outputs variables Zhang and  
 218 Rundell (2006). In addition to first order effects, the variance decomposition methods, also allows the  
 219 quantification of second order effects sometimes referred as total order effects. Total order effects indices  
 220 are useful for the assessment of the interaction between factors which cannot be expressed by a simple  
 221 lineal superposition.



**Figure 2. The different types of sensitivity analysis and their associated methodologies and techniques.**

222 One of main drawbacks for applying variance decomposition methods on large spatially explicit  
 223 individual-based models is the requirement of very high number of model evaluations in order to produce  
 224 consistent results Herman et al. (2013). An alternative approach, in those cases where it is impractical

225 or computationally unfeasible a fully quantitative analysis, is the application of the Morris screening  
 226 method. The Morris method deliver qualitative information allowing to rank the importance of input  
 227 factors requiring lees model evaluations, which in some case can one order of magnitude inferior to the  
 228 Sobol method Saltelli (2008).

229 The Sobol is a method for sensitivity analysis based on the decomposition of the variance of model  
 230 output and is particularly suitable for discovering the effect of high order interactions between input  
 231 factors. The interaction means non-linearity where the total effect of two input factors  $x_1$  and  $x_2$  on the  
 232 model output  $Y$  are not equivalent to the sum of the individual effects. The general form of sensitivity  
 233 indices for Sobol methods are shown in Equation (1) and Equation (2), respectively the first order and  
 234 total order indices.

$$S_i = \frac{V_i}{V(Y)} \quad (1)$$

$$S_{Ti} = 1 - \frac{V(Y) - V_i}{V(Y)}, \quad (2)$$

235 where the terms  $V_i$  and  $V(Y)$  are respectively the variance contribution attributed to the  $i$ th parameter and  
 236 the total variance. The expression  $V(Y) - V_i$  represents the total variance with exception of the variance  
 237 which is generated by the parameter  $i$ . The total order index  $S_{Ti}$  is the contribution of all input parameters  
 238 but one, the  $i$ th parameter, and hence estimating the effect of that parameter on the variance reduction  
 239 Saltelli (2008).

240 The total variance  $V(Y)$  for a model with  $n$  input parameters can be expressed as shown in Equation  
 241 (3) as long as the orthogonality of input factors precondition holds.

$$V(Y) = \sum_i V_i + \sum_{i<j} V_{ij} + \sum_{i<j<k} V_{ijk} + \dots + V_{12\dots n}, \quad (3)$$

242 being  $V(Y)$  the total variance from model output and the components  $V_i$ ,  $V_{ij}$  and  $V_{ijk}$  respectively the  
 243 variance contribution from the parameter  $i$ , the variance contribution form input parameters  $i$  and  $j$  and  
 244 the variance contribution form input parameters  $i$ ,  $j$  and  $k$ . Finally, the component  $V_{12\dots n}$  express the  
 245 interactions from all parameters present in the model.

246 The application of Sobol method, as have been mentioned, can be computationally expensive and  
 247 sometimes could be useful to reduce the problem dimensionality filtering only the most significant  
 248 parameters or even simplifying the model structure considering only the parameters accounting for the  
 249 most of the variability in the model output. It can be accomplished using the Morris screening method to  
 250 rank the importance of input parameters. The Morris method is an OAT method, meaning that it changes  
 251 just one factor keeping all other input parameters fixed. The input factors are allowed to vary in discrete  
 252 levels within the relevant parameter range Morris (1991). The method is considered to be more effective  
 253 when the number of most significant input parameters are a small subset of model parameters Saltelli et al.  
 254 (2004).

255 The original work of Morris Morris (1991) define two metrics for ranking input factors which are  
 256 depicted by  $\mu$  and  $\sigma$  values<sup>1</sup>. Further, another metric termed  $\mu^*$  has been suggested by Campolongo et al.  
 257 (2007) which use absolute values in order to handle effects of distinct signs canceling each other. These  
 258 metrics for ranking input factors are calculated from what has been termed elementary effects. Therefore,  
 259 considering a model with  $k$  input parameters and being  $x = (x_1, x_2, \dots, x_k)$  any value from the region of  
 260 experimentation  $\Omega$ , the elementary effects are calculated according to the Equation (4).

$$ee_i(x) = \frac{y(x_1, \dots, x_{i-1}, x_i + \Delta, x_{i+1}, \dots, x_k) - y(x)}{\Delta} \quad (4)$$

261 The region of experimentation  $\Omega$  is a grid defined by the number of  $k$  input factors and by the  $p$   
 262 discrete levels for every parameter. The recommendations for the values of  $p$  and  $\Delta$  are respectively that  
 263 the first should be an even number of levels and the second calculated by the expression  $\Delta = p/(2(p-1))$

<sup>1</sup>Not to be confused with population mean and standard deviation

264 Morris (1991) Saltelli et al. (2004). The value of  $\Delta$  has important implications in the model analysis. It  
 265 has been shown that in some situations choosing an alternative value calculated as  $\Delta = 1/(p - 1)$  can  
 266 detect non-monotonic behaviors that the suggested standard calculation are not able to capture otherwise  
 267 van Houwelingen et al. (2011).

268 The metrics of Morris method are calculated over the  $F_i$  and  $G_i$  distributions for every input parameter.  
 269 These distributions are generated taking random samples of  $x$  from  $\Omega$  for calculating the elementary  
 270 effects and the only difference between them is that  $G_i$  uses the absolute values of elementary effects  
 271  $|ee_i(x)|$  as described in Campolongo et al. (2007) Saltelli (2008). The estimation of Morris metrics are  
 272 carried out by taking  $r$  samples from  $F_i$  and  $G_i$  distributions according to the Equations (5), (6) and (7).

$$\mu = \sum_{i=1}^r \frac{ee_i(x)}{r} \quad (5)$$

$$\mu^* = \sum_{i=1}^r \frac{|ee_i(x)|}{r} \quad (6)$$

$$\sigma = \sqrt{\sum_{i=1}^r \frac{(ee_i(x) - \mu)^2}{r}} \quad (7)$$

273 These three metrics can be used to extract valuable information about the model behavior, in addition  
 274 to ranking the input factors. For instance, a low value of  $\mu$  and a high value of  $\mu^*$  is high, points that the  
 275 input factor under scrutiny, possibly has a non-linear behavior having different signs in function of the  
 276 system trajectory Saltelli et al. (2004). A high value of  $\mu$  indicates that the input has a monotonic effect  
 277 on the model output.

278 The sensitivity analysis methods require significant samples from input space in order to provide  
 279 reliable results. It is customary to choose between some experimental design Hicks (1993) for generating  
 280 the collection of input parameters needed by evaluating the model and allocating the variance contribution  
 281 of every model parameter. The most generally applied sampling schemas are based on random sampling,  
 282 full factorial designs or Latin hypercube sampling.

## 283 OVERVIEW OF R/REPAST PACKAGE

284 In the previous sections we had seen some fundamental ideas on model building and the role occupied  
 285 by sensitivity analysis methods in the iterative modeling life-cycle. We have also introduced the basic  
 286 principles of sensitivity analysis focusing on two main techniques namely the Morris Elementary Ef-  
 287 fects Screening Morris (1991) and the Sobol GSA method for variance decomposition Saltelli (2008).  
 288 Both methods have a wide range of applicability, making them suitable for their use in the analysis of  
 289 Individual-based models. These methods require the model be evaluated many times with a different set  
 290 of input parameters, making completely impractical undertaking a manual analysis introducing individual  
 291 parameters manually on a graphical user interface. The Repast is an extremely flexible framework for  
 292 object-oriented development of Agent-based models using Java language but it lacks from model analysis  
 293 tools. On the other hand, the GNU R is a superb open source tool for data analysis with a vast and  
 294 active community developing and adding new methods to the core R system. Bearing this in mind, we  
 295 introduce our package R/Repast which bring together the best of both worlds. Roughly speaking, the  
 296 R/Repast package have two main objectives: (a) Provide an interface for running Repast models from R  
 297 and gathering the simulation data generated and (b) Automating the application of sensitivity analysis and  
 298 simple model calibration methods to the Repast models. The R/Repast is an open source project delivered  
 299 under the MIT license system. The package provides a powerful and simple R API<sup>2</sup> which reduces the  
 300 code required for running the most commonly used experimental methods suitable for . The software  
 301 and the user manual can be downloaded from CRAN website and the complete project source code from  
 302 GitHub repository. Both are available respectively from the following URLs:

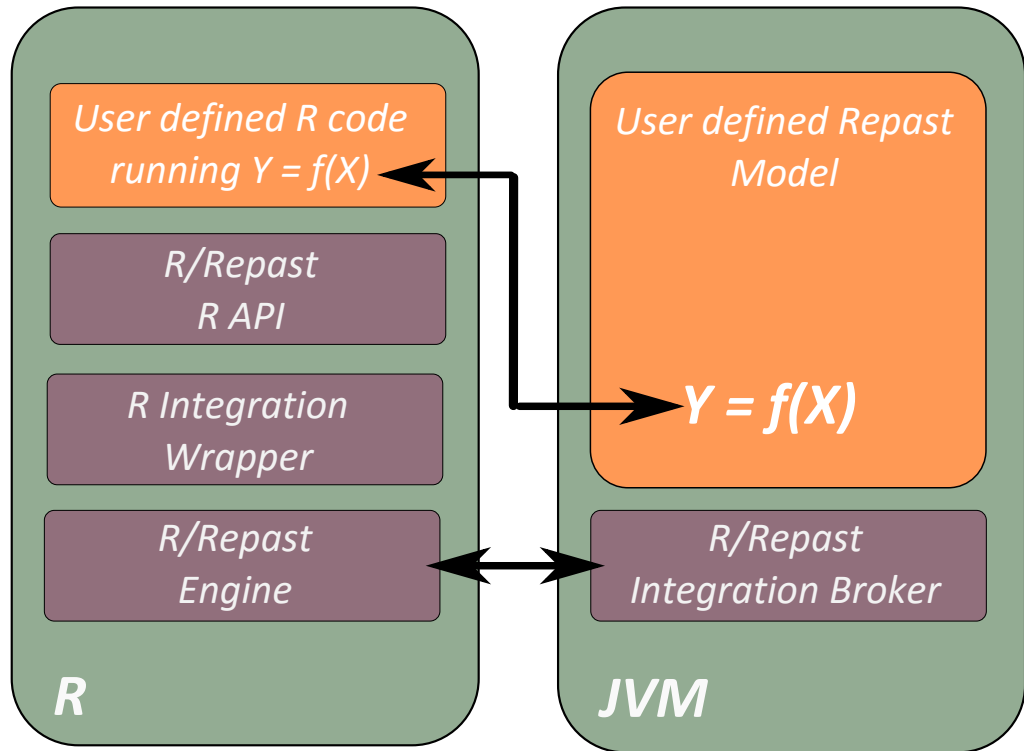
- 303 • <https://cran.r-project.org/web/packages/rrepast/>
- 304 • <https://github.com/antonio-pgarcia/RRepast>

<sup>2</sup>Application Programming Interface



305 **Design**

306 The **R/Repast** was intended primarily for invoking Repast Simphony models from inside GNU R  
307 environment. Additionally, the package contains more high level and value added features for experimental  
308 design and experiment analysis to address the specific need of individual-based models. The underlying  
309 implementation idea is to provide a set of turnkey features for facilitating the task of applying the  
310 sensitivity analysis to models. Functionally, the package consists in four modules which interoperate  
311 together for instantiate and running the Repast code inside R. These four components are (a) the Repast  
312 Integration Broker, (b) the Repast Integration Engine, (c) The R Integration wrapper and finally, (c) the R  
313 API for Experiment design. An schematic view of package architecture is shown in Figure 3.



**Figure 3.** The R/Repast general architecture. The scheme shows in the left box the R environment and the associated components of R/Repast. The right box represents the Repast Simphony model running within a Java Virtual Machine as well as the R/Repast integration broker component.

314 The R/Repast integration broker and the R/Repast engine are both written in java code and are required  
315 for instantiating and loading the Repast Simphony model in batch mode. The R/Repast engine contains  
316 also the required hooks for transferring the model output data from Java to R environment. The engine can  
317 transfer data from aggregated data set defined by the modeler on the Repast model. An aggregated data  
318 set is a Repast Simphony entity used for collect data about the simulation model agents which can be used  
319 for plotting or saving the model output data to a file using a file sink. A *File Sink* is Repast component for  
320 saving simulation data to a file. The aggregated datasets use some kind of aggregate operation, such as  
321 counting, averaging, summing or any other used defined aggregate operation North et al. (2013c), North  
322 et al. (2013a), North et al. (2013d). The R integration wrapper is the R code for linking together the R  
323 and Repast subsystems. This module consist in several wrapping functions for encapsulating the Java  
324 code calls implemented using the *rJava* package Urbanek (2016). These functions are prefixed with the  
325 [Engine] keyword and, although exported in the R/Repast package, they are not intended for general use.

326 **The R/Repast R API**

327 The module entitled **R/Repast R API** is the primary entry point for the user defined code and relies on  
328 the subsystems mentioned previously for providing three group of functionalities for facilitating modelers  
329 to analyze the simulation output. These group functionalities are the following:

- 330 • Execution and control of Repast Symphony code.
- 331 • Basic functions for experimental design.
- 332 • High level functions for a complete experiment in one call.

333 The functionalities on the first group are those required for the basic interface between Repast and  
 334 R system, such as instantiating and running a Repast Symphony model, retrieving the declared model  
 335 parameters, getting their default values, setting parameter values as well as running basic experimental  
 336 designs and saving simulation data. The list of these functions are shown in Table 1.

337

**Table 1. The basic R/Repast API Functions. These functions are used for loading and modifying the default parameters defined for model and also for running the simulation.**

Function name	Description
<b>Model(d, t, o, l)</b>	This function creates an object instance for linking the Repast model to an R object. The required parameters are the directory where the model has been installed ( <i>d</i> ), the duration of simulation in Repast ticks ( <i>t</i> ), the name of any aggregated dataset of model for draining data generated by the model simulation ( <i>o</i> ) and a Boolean flag ( <i>l</i> ) which tells the function to call the Load method. The default value is FALSE.
<b>Load(m)</b>	This function loads the Repast scenario from model's directory. The only required parameter ( <i>m</i> ) is an instance of Repast Model created with previous function.
<b>Run(m, r, s)</b>	The purpose of this function is to execute a single round of simulation using just one parameter set. The parameters for this function are a model instance ( <i>m</i> ), the number of repetitions ( <i>r</i> ) and a collection of random seeds ( <i>s</i> ) to be used for each one of the repetitions. The only required parameter is the model instance, created with the <i>Model()</i> function. The default value for <i>r</i> is one.
<b>RunExperiment(m, r, d, F)</b>	Execute a complete experimental setup for different sets of parameters. The parameters required are a model instance ( <i>m</i> ), the number of replications ( <i>r</i> ), the experimental design ( <i>d</i> ) and finally a user provided calibration function ( <i>F</i> ). The experimental design parameter is an <b>R</b> data frame containing a complete set of model's parameter per row. The function returns a list with three data frame elements: the <i>paramset</i> , the <i>output</i> and <i>dataset</i> which holds respectively all simulated input parameters, the result of user provide calibration function and the complete dataset produced during the experiment execution.
<b>GetSimulationParameters(m)</b>	Returns the complete list of parameters declared by the model. The parameter ( <i>m</i> ) is an instance of Repast model generated with <i>Model()</i> call described previously.
<b>SetSimulationParameters(m, p)</b>	Modify several parameters at once.
<b>SaveSimulationData(t, e)</b>	Exports the results of Run or RunExperiment to a csv or excel files. The parameters <i>t</i> and <i>e</i> are respectively the format of exported data (xls or csv) and the experiment results returned by <i>RunExperiment()</i>

338

339 The second group of methods within R/Repast R API contains the functionalities required for setting  
 340 up and applying a complete experimental design to a Repast simulation model. The group include  
 341 functions for adding the input factors and the relevant input range which the modeler wants to evaluate.  
 342 The group also have functions for generating the experiment inputs using different sampling approaches.  
 343 It is not required to add as input factors all declared model parameters, the modeler can just evaluate a  
 344 small subset keeping the other factors fixed. The functions of this group are presented in Table 2.

345

**Table 2. The Experimental Setup API functions. These functions are used for experimental design, parameter calibration and sensitivity analysis.**

Function name	Description
<b>AddFactor(f, l, k, b, u)</b>	Creates the parameter collection for the experimental setup. The function requires the data frame ( <i>f</i> ) where parameter will be added, if this parameter is not provided a new data frame will be created. The second parameter ( <i>l</i> ) is the random function used internally, the default value is <i>runif</i> which will be the valid choice in many cases, the next parameter ( <i>k</i> ) is the name of factor, the value provided must match some parameter defined in the repast model. The following two parameters ( <i>b</i> ), ( <i>u</i> ) are the lower and the upper range, respectively. The function returns the updated ( <i>f</i> ) data frame with the new parameter.
<b>AoE.RandomSampling(n, f)</b>	Also known as Monte Carlo sampling, generate an experimental design based on making random samplings of parameter space. The function takes two parameters, the sample size ( <i>n</i> ) and the factor ( <i>f</i> ) data frame created using <b>AddFactor()</b> . The function returns the design matrix for the provided parameters.
<b>AoE.LatinHypercube(n, f)</b>	Generates an experimental design using the Latin Hypercube stratified sampling technique which is a more efficient sampling scheme, in terms of model evaluations, than the pure random sampling. The parameters ( <i>n, f</i> ) and return values are the same already described for the function <b>AoE.RandomSampling()</b> .
<b>AoE.FullFactorial(n, f)</b>	Creates a factorial design where the effects of all independent variables of model are studied simultaneously, which implies many more model evaluations. The parameters ( <i>n, f</i> ) and return values are the same already described for the function <b>AoE.RandomSampling()</b> .
<b>BuildParameterSet(d, p)</b>	Constructs the data frame required for executing <b>RunExperiment()</b> . The function takes two parameters: the design matrix ( <i>d</i> ) created with one of previous functions and the declared parameters ( <i>p</i> ) defined in the Repast Model with the default values retrieved using the function <b>GetSimulationParameters()</b> . The functions returns a data frame with varying and fixed parameters for the experimental setup of choice.

346

347 Finally, the third group contains the "Easy" API functions. These functions are intended to provide  
 348 a complete method implementation which is accessible using just one R function call. The user has  
 349 to provide the directory location where the Repast model is installed, the objective function and the  
 350 parameters relevant to the specific method. The currently available Easy API methods are presented in  
 351 Table 3. The objective function is a user defined R function over the model output for calculating and  
 352 returning a cost metric for the simulation outputs of interest. The return of objective functions is the target  
 353 for the application of the analysis method.

**Table 3.** The easy API functions. These functions are the preferred entry point for the eventual users. These "Easy" functions lump together a complete experiment task in just one call, reducing the number of lines of code required.

Function name	Description
<b>Easy.Stability(d,o,t,f,s,r,v,F)</b>	Evaluate the behavior of model output in order to determine the minimum required number of replication of the chosen experimental setup. The function accept the following parameters: the model installation directory ( <i>d</i> ), the aggregated data source defined within the Repast model ( <i>o</i> ), the simulation time in Repast ticks ( <i>t</i> ) which default value is 300 ticks, the input factors to be sampled ( <i>f</i> ) created with the previously mentioned function <b>AddFactor()</b> , the number of parameter samples ( <i>s</i> ), the desired number of replications to be tried ( <i>r</i> ) being the default value 100, the output variables of interest which will be checked for their stability and convergence of the coefficient of variation ( <i>v</i> ); if this parameter is leaved empty all output variables are checked and finally the user provided calibration function ( <i>F</i> ) for determining the best input parameter combination.
<b>Easy.Morris(d,o,t,f,p,s,r,F)</b>	This function performs all required tasks for carrying out the method of Morris for screening. The parameters are practically the same as described for the previous function with exception of parameters ( <i>p</i> ) and ( <i>s</i> ) which are respectively the levels of input factors and the number of sampling points of Morris method Pujol et al. (2015).
<b>Easy.Sobol(d,o,t,f,n,r,F)</b>	Encapsulate all required steps for performing sensitivity analysis using Sobol method. The method of Sobol is a global sensitivity analysis technique based on the decomposition of output variance (Saltelli et al. (2004); Pujol et al. (2015)). The parameter semantics are the same already described: the model installation directory ( <i>d</i> ), the aggregated data source defined within the Repast model ( <i>o</i> ), the simulation time in Repast ticks ( <i>t</i> ), the input factors to be sampled ( <i>f</i> ), the sample size ( <i>n</i> ), the desired number of replications ( <i>r</i> ) and calibration function ( <i>F</i> ).
<b>Easy.Calibration(d,o,t,f,n,r,F)</b>	This function estimates the best set of input parameters or factors, performing a set of model executions in order to sample the calibration function. The objective of this function is to minimize the output of calibration function provided by the user.
<b>Easy.Setup(d,l)</b>	The parameters ( <i>d</i> ) and ( <i>l</i> ) are respectively the directory where repast model is installed and the location of R/Repast deployment directory. If omitted, it assumes as the default value, the directory where the Repast model is installed. The function is required for automatically making the changes in the model configuration for adding the integration code, for deploying the Java <i>jar</i> files with the integration code and for preparing the deployment directory. That directory will hold the JVM logs and the saved model output data sets.

### 356 The objective function interface

357 The last piece of R/Repast architecture is the definition of the objective function which actually allows the  
 358 flexible definition of the model analysis target decoupling it from the Repast dataset output. As we have  
 359 mentioned previously, any model is a functional relationship between a vector of input parameters *X* and

360 a scalar dependent variable  $y$  and expressed as  $y = f(X)$ . On the other hand, usually the dataset collected  
361 from Repast model execution will be a time series where the aggregated measure will be collected at fixed  
362 intervals. Therefore, some transformation must be applied in order to obtain a value consistent with the  
363 functional definition. In addition, even though the value returned from the Repast model were a scalar  
364 one, it would add much more maintainable and flexible a to act upon it directly from R without have  
365 to making changes in the Repast code. The objective function is also necessary for calibrating, where  
366 the output values are compared to some reference data or even for more complex tasks, such as tuning  
367 oscillations in the population output. It is also the place for normalizing the model outputs. The objective  
368 function is a required parameter for all methods presented here.

369 The specification of R/Repast requires the objective function having two input parameters. The  
370 first input parameter for the objective function is the input parameter set used for executing the Repast  
371 model, the second parameter are the results generated by executing the model and corresponding to and  
372 aggregated data set in the Repast model. The objective function must return one or more scalar values  
373 grouped using the `cbind()` Crawley (2007) R function. The complete function signature is shown in  
374 Figure 4.

```
1 objective<- function(params , results) {  
2   cbind( ... )  
3 }
```

**Figure 4. The skeleton of objective function. The function has two parameters and must return a one or more scalar values.**

## 375 EXAMPLES OVERVIEW

376 In the next sections we will provide examples on how the R/Repast can help modelers on the analysis of  
377 their simulation models. Three examples will be used for illustrating the application of some package's  
378 functionalities and what kind of information these functions can offer about the simulation outputs. For  
379 clarifying what every model does a summary version of ODD will be given for facilitating a general idea  
380 about these models. The Overview, Design concepts and Detail (ODD) is a protocol Grimm et al. (2006)  
381 Grimm et al. (2010) which has been proposed as a standard way to specify and describe Individual-based  
382 models. A brief description on the model structure and parameters will be given in order to allow the  
383 readers to understand the kind of questions the model is intended to answer and how R/Repast is can be  
384 used for analyzing the model outputs. The last section for each model under the title of **Model analysis**  
385 is not part of ODD protocol but it is included to show the results of running the R/Repast model analysis  
386 methods.

387 The first model used as example here is a spatially explicit individual-based representation of bacterial  
388 conjugation using BactoSIM for simulating the plasmid spread on a surface attached bacterial colony  
389 Prestes García and Rodríguez-Patón (2015). The example will be used for showing the application  
390 of *Easy.Stability* method for finding the number of replications of simulation experiments required for  
391 obtaining consistent outputs. The second example is a Repast implementation of the omnipresent predator-  
392 prey model describing the interaction between two species. This one is part of examples coming along the  
393 standard Repast distribution and will be used for showing the application of *Easy.Morris* function. Finally,  
394 the third example is an instance of the common pool problem in the context of bacterial conjugation.  
395 This model was developed exclusively for this work. This model will be used for exemplifying the use  
396 of *Easy.Sobol* method. The complete sources for all projects are available respectively in the following  
397 locations:

- 398 • **BactoSIM:** <https://github.com/antonio-pgarcia/haldane>
- 399 • **Predator-Prey:** The sources come with the Repast distribution.
- 400 • **T4SS Common Pool:** <https://github.com/antonio-pgarcia/PoolT4SS>

401 For convenience, in order to facilitate the experiments shown in this paper being reproduced elsewhere,  
402 we also provide the pre-built installers for the three projects mentioned previously. The installers can be  
403 download from URL shown below:

- 404 • **BactoSIM:** <http://goo.gl/YYIt1o>
- 405 • **Predator-Prey:** <http://goo.gl/cJ5z2r>
- 406 • **T4SS Common Pool:** <http://goo.gl/zq4LH0>

407 In order to reproduce the examples shown in the next sections, it is required a computer with a  
408 Java JVM and GNU R installed. The examples have been produced and tested on a windows box with  
409 java 1.8 and GNU R 3.3.1. If these preconditions are met, just proceed to download and install the  
410 examples and the R/Repast package. The installation of R/Repast is carried out using the install command  
411 `install.packages("rrepast")` on the R environment. Once the previous steps have been completed, just  
412 copy and paste the examples shown in this paper, taking care of changing the references to the model  
413 installation directory to the directories where the models have been installed locally.

## 414 **EXAMPLE 1: BACTOSIM**

415 Normally, one of the advantages of using individual-based models for biological or ecological processes  
416 is the possibility of incorporating variability at an individual level. Therefore, unlike deterministic model,  
417 in order to get trustworthy results, the simulation must be repeated a number  $N$  of times to achieve stable  
418 value on the output variance. The objective of the first example is to show the application of a simple  
419 method for finding the minimal number of replications of a simulation model which is required for the  
420 variance of response variables become stable, converging to a common value. A straightforward way  
421 to determine the output stability has been suggested in Thiele et al. (2014) Lorscheid et al. (2012) and  
422 consists in to compute the coefficient of variation<sup>3</sup> of the output of interest with and increasing number  
423 of repetitions while keeping the input parameters fixed. The number of replications for which the values  
424 of coefficient of variation stop to vary are the minimum number of repetitions necessary for getting robust  
425 results. In R/Repast we have implemented that method which is accessible through the *Easy.Stability* API  
426 call.

427 For this example, the BactoSIM Prestes García and Rodríguez-Patón (2015) model will be used. This  
428 is an individual-based model of bacterial conjugation process. The bacterial conjugation is a form of  
429 lateral genetic transfer which occur naturally in bacterial colonies Arutyunov and Frost (2013). The  
430 conjugation consists in the transference of a conjugative plasmid from a donor cell to a recipient cell. The  
431 plasmids are small circular DNA sequences which replicates independently from the main chromosome  
432 of their hosts Bergstrom et al. (2000). The conjugation is considered one the causes of the rapid evolution  
433 and adaptation of bacterial colonies and the spread of antibiotic resistance Chen et al. (2005) Slater et al.  
434 (2008). The BactoSIM model is currently being used for an evaluation of the main factors governing  
435 the plasmid dispersion. A preliminary evaluation has shown that the point in the cell cycle are the  
436 principal factor responsible for the global dynamics of plasmid infective dispersion Prestes García and  
437 Rodríguez-Patón (2015) which is consistent with some observations Seoane et al. (2011) taken from  
438 individual bacterial cells.

### 439 **Model description**

440 The model description follows the ODD (Overview, Design concepts and Detail) protocol for describing  
441 individual-based models Grimm et al. (2006) Grimm et al. (2010). The model is implemented in java  
442 language using Repast Symphony agent-based simulation framework North et al. (2013b).

#### 443 ***Purpose***

444 The objective of this model is the assessment of the best strategy for modeling and implementing the  
445 conjugation rule which provides the best fit to experimental data and better captures the most plausible  
446 process structure.

#### 447 ***Entities, State variables and scales***

448 The model comprises two entity types, namely the bacterial individuals or agents and environment. The  
449 environment contains the rate limiting number of nutrient particles required for the cell metabolism  
450 and growth. All agents evolve in a computational domain defined by a  $1000 \times 1000 \mu m$  squared lattice

---

<sup>3</sup>Also known as relative standard deviation given by  $CV = \sigma/\mu$  which provides a normalized version of the standard deviation expressed relatively to the output mean

451 divided in  $10^6$  cells of  $1 \times 1 \mu\text{m}$  representing a real surface of  $1\text{mm}^2$ . In this model the *agents* representing  
 452 bacterial cells are defined individually by two main state variables, namely the *plasmid infection state*  
 453 and the  $t_0$ . The plasmid infection states are  $\mathcal{Q} = R, D, T$  and the respective transition function for  
 454 **conjugative** plasmids,  $\delta$  is shown in (8). For the *oriT* construction only the first transition rule applies  
 455 since transconjugant cells are sterile. The  $t_0$  is the time of cell birth or the time of the last cellular division,  
 456 it is employed in the estimation of agent doubling time used in the division decision rule. The T4SS pili  
 457 is also taken into account and the agents have a state variable representing the number of pilus already  
 458 expressed and available in cell surface.

$$\delta = \begin{cases} (D, R) \rightarrow (D, T) \\ (T, R) \rightarrow (T, T) \end{cases} \quad (8)$$

459 Finally, the environment will hold the initial nutrient concentration for every lattice cell. In the model  
 460 initialization, a fixed amount of substrate particles will be distributed evenly over all lattice sites.

#### 461 **Process overview and scheduling**

462 The dynamics of bacterial conjugation is modeled as the execution of following set of cellular processes:  
 463 the cellular division, the T4SS pili expression, the shoving relaxing which avoid bacterial cells to overlap  
 464 and allow a more realistic colony growth and the conjugation process. The state variable update is  
 465 asynchronous. The order of execution of this process is shuffled to avoid any bias due to a purely  
 466 sequential execution of model rule base, see 5. The conjugation process is modeled in three different  
 467 ways with respect to the time when conjugation event is most prone to happen and the results are  
 468 compared. Thus the conjugation is defined by to variables: the value of intrinsic conjugation rate ( $\gamma_0$ )  
 469 which determines how many transfers should be performed by a single bacterial cell and the cell cycle  
 470 point which defines the time when the conjugative events are likely to occur.

471 The model input and initialization requires the parameters shown in Table 4. The *costT4SS* is the  
 472 total cost of pili expression. The cost applied for a single pilus expression is *costT4SS/param(maxpili)*.  
 473 The *param(maxpili)* is actually a constant having the value of 5 for *E. coli*. The *cellCycle* parameter  
 474 indicating two things: the type of modeling rule and its parameter. A value of  $-1$  set the model to  
 475 conjugate as soon an infected cell finds a susceptible one. Setting the parameter to 0 will randomize  
 476 de conjugation time between  $t_0$  and  $G$ . Finally using a value greater than zero indicates the specific  
 477 point in the cell cycle for conjugation. A polynomial equation fitted to the experimental data where the  
 478 dependent variable represents the conjugation rate  $T/(T + R)$ . Setting *isConjugative* flag to false creates  
 479 a simulation where the transconjugant cells are sterile, in other words they are unable to conjugate. The  
 480 *equation* is used only for comparing the quality of simulation output.

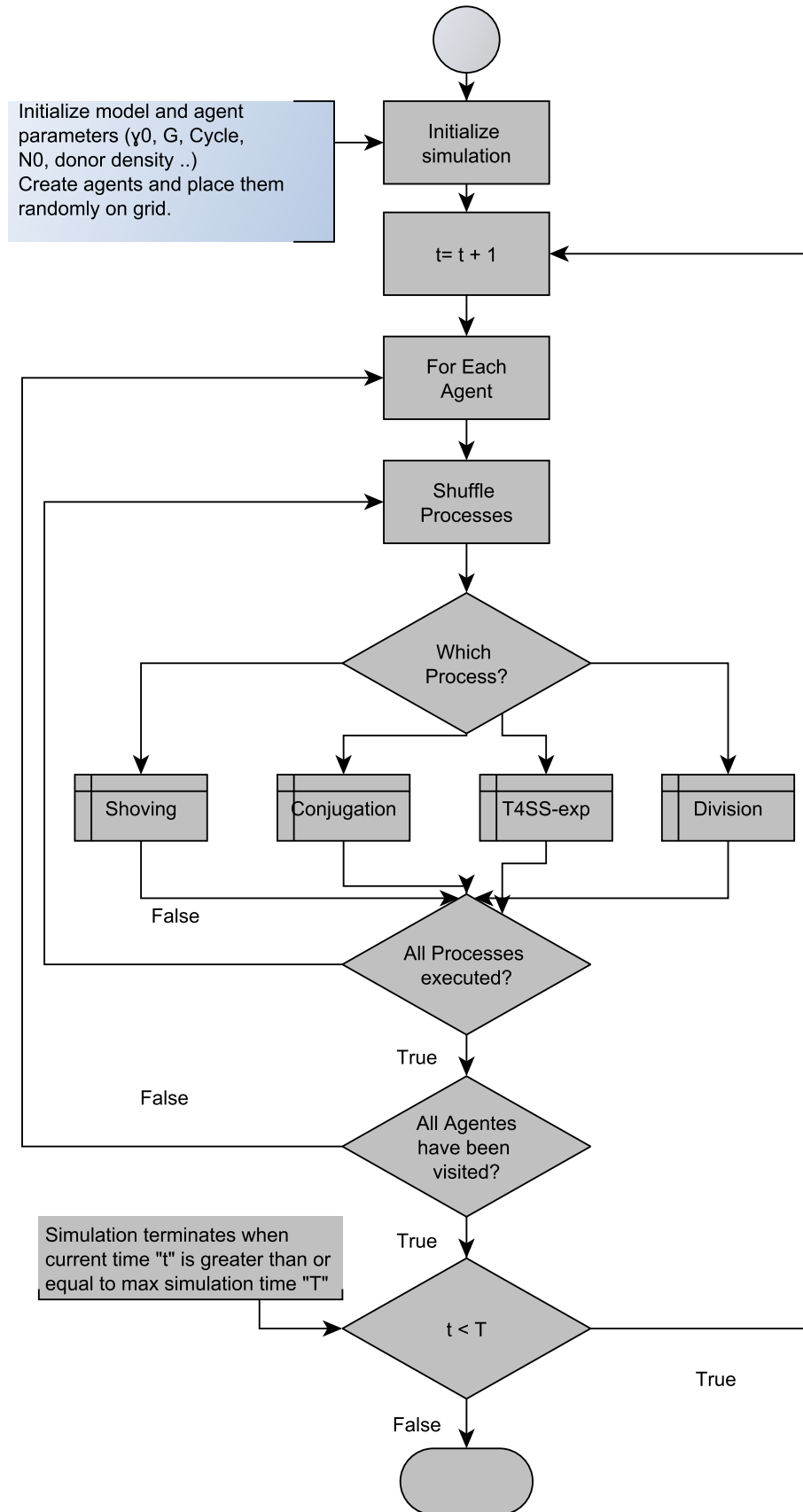
481

**Table 4. The complete list of model initialization parameters.**

Parameter	Unit	Description
G	minutes	Average doubling time for plasmid free cells
cellCycle	% of G	The percentage of cellular cycle for conjugation
costConjugation	% of G	The penalization due to a conjugative event
costT4SS	% of G	The Pilus expression cost
$\gamma_0$	conjugations/cell	Upper limit for conjugations performed by an agent
isConjugative	true—false	Defines a conjugative or a mobilizable plasmid
isRepressed	true—false	The T4SS expression state for the plasmid
$N_0$	cells/ml	Initial population expressed in cells/ml
donorRatio	% of $N_0$	The initial density of donor cells ( $D$ )
Equation	N/A	An equation for experimental data

#### 483 **Design concepts**

484 **Basic Principles** — Three models differing in the way the conjugation rule is implemented and their  
 485 results compared to the available experimental data. The best strategy can be used to build models which



**Figure 5.** The flow diagram showing the overview on how bacterial process are scheduled in the BactoSIM simulation model.



486 could serve as a predictive tool for synthetic biology and to explore some aspects which are hard to  
487 observe directly in experimental studies of plasmid spread. The key points of this model lies on the idea  
488 of the existence of a local or intrinsic conjugation rate, which has been termed  $\gamma_0$ . This intrinsic rate  
489 stands for the number of plasmid transfer events, or conjugations on a cell life-cycle basis. In addition,  
490 the global infective wave speed depends directly from the specific point in the bacterial cell cycle when  
491 conjugative event is triggered.

492 **Emergence** — The model intends to find out what will be the global outcome arising as function of local  
493 rules defining the evolution of the bacterial cells and their interaction with adjacent neighbors. With this  
494 objective, the model incorporates the most significant aspects of the spatial structure and the behavior  
495 of the cellular processes that are related to the conjugation. Specifically, the values of the generation  
496 time of donor and transconjugant cells are one of the emergent properties depending from the metabolic  
497 penalizations applied both for conjugation event and for the expression of T4SS genes.

498 **Adaptation** — All agents adapt their growth according to the local availability of nutrient and space.

499 **Fitness** — It is considered implicitly to the extent that plasmid free individuals will present a better  
500 adaptation in terms of growth rate than plasmid bearing cells.

501 **Prediction** — The model is intended to provide prediction regarding the range of possible values for the  
502 number of plasmid transfer events per cell cycle and the cell cycle point when conjugative transfer is most  
503 likely to happen.

504 **Sensing** — All process defined over the agents implicitly sense the local environment and the close  
505 neighborhood for their decisions.

506 **Interaction** — Bacterial cells interact with their nearby individuals for nutrient access, cellular division,  
507 mate pair formation and plasmid transfer.

508 **Stochasticity** — Stochasticity is introduced at individual level for all cellular process sampling a normal  
509 deviate and fitting the value to corresponding process.

510 **Collectives** — No collectives are taken into account in this model.

511 **Observation** — The output target variables will be saved at intervals of one minute of simulated time.

512

### 513 **Initialization**

514 The simulation model is initialized with a population of plasmid free ( $R$ ) and plasmid bearing ( $D$ ) cells  
515 according to input parameters. The agents are placed randomly within a circular surface centered over the  
516 lattice central position. The radius of circle where agents are placed is calculated as function of  $N_0$  in  
517 order to be consistent to the desired initial cell density Zhong et al. (2012). The simulation environment is  
518 also initialized with a number of nutrient particles in order to support the half of the estimated number of  
519 cellular divisions and the rationale behind it is to capture the intercellular competition for nutrient access.

### 520 **Model analysis**

521 The objective of stability analysis is to find the minimum number of experimental setup replications  
522 required for achieving reliable results. Thereby, the model output response is evaluated for an increasing  
523 number of repetitions allowing the evaluation of the convergence for output variance of simulation outputs.  
524 The complete listing for carrying out the stability check for the BactoSIM model is shown in Figure 6. As  
525 can be observed, the complete implementation of model analysis encompasses five steps. These steps  
526 are conserved for all high level functions available in R/Repast package. The **step 0** clean all existing R  
527 objects, loads the R/repast package and set the random seed for the analysis. The **step 1** is the definition of  
528 the objective function which can be any user provided function following the R/Repast API specification.  
529 It is not strictly necessary for the **Easy.Stability** as the coefficient of variation is calculated for the model  
530 output variables. In this example the objective function is basically the comparison of simulated data  
531 and experimental data using the normalized root mean square error API call *AoE.NRMSD*. The **step 2**  
532 is adds the model input factors for which the importance on the model output will be assessed and their  
533 biologically relevant range of variation. It is necessary to add at least one parameter which will be varied,  
534 while all other model parameters are keep fixed using the default value or with a value previously set using  
535 the R/Repast API *SetSimulationParameter*. The purpose of **step 3** is to configure automatically the Repast  
536 model with the integration broker and for initializing the integration directory. Finally, the **step 4** is where  
537 the analysis method is invoked, all analysis methods will return a list holding three objects, namely the  
538 *experiment*, the *object* and the *charts*. The experiment contains simulation parameters and results, the

539 object is method specific and finally the charts are pre-generated graphs for the method results<sup>4</sup>

```
1 # Step 0
2 rm(list=ls())
3 library(rrepast)
4 set.seed(161803398)
5
6 # Step 1
7 objective<- function(params, results) {
8   Rate<- AoE.NRMSD(results$X.Simulated, results$X.Experimental)
9   cbind(Rate)
10 }
11
12 # Step 2
13 f<- AddFactor(name="cyclePoint", min=40, max=90)
14 f<- AddFactor(factors=f, name="conjugationCost", min=0, max=30)
15 f<- AddFactor(factors=f, name="pilusExpressionCost", min=0, max=30)
16 f<- AddFactor(factors=f, name="gamma0", min=1, max=4)
17
18 # Step 3
19 Easy.Setup("/models/BactoSim(HaldaneEngine-1.0)")
20
21 # Step 4
22 r<- Easy.Stability("/models/BactoSim(HaldaneEngine-1.0)", "ds::Output", 300, f, 1, 100, c(), objective)
```

**Figure 6.** The listing for stability of model output method using the *Easy.Stability* function from R/Repast.

540 The method will generate automatically one chart for each model output<sup>5</sup>. One of the output chart  
541 of model is shown in Figure 7 for the variable named *X.Simulated*. As can be observed, the coefficient  
542 of variation of these variable decreases as the sample size increases. The variation starts to become  
543 acceptable with a sample size of 25 and approximately with sample size of 50 we can see that coefficient  
544 of variation become stable. Therefore, we can feel relatively confident with or model results with a  
545 number of replications greater than 25. Of course it is important to take into account the computation cost  
546 of our model in order to select a value for the number of repetitions.

## 547 EXAMPLE 2: PREDATOR-PREY

### 548 Model description

#### 549 Purpose

550 The purpose of Predator-Prey model presented here is to provide an alternative individual based-model  
551 implementation for the classic ODE model describing the association between two species. The model  
552 will be used to show the application of Morris method for ranking the most important parameters.

### 553 Entities, State variables and scales

554 The model comprises three entities or agent types, the wolves, the sheep individuals and the grass. These  
555 agents evolve in a computational domain of a  $50 \times 50$  units with periodic boundaries, representing a large  
556 portion of space. The agents are positioned in a continuous bi-dimensional space and are free to move.  
557 On the other hand, the grass agent is placed in a discrete grid.

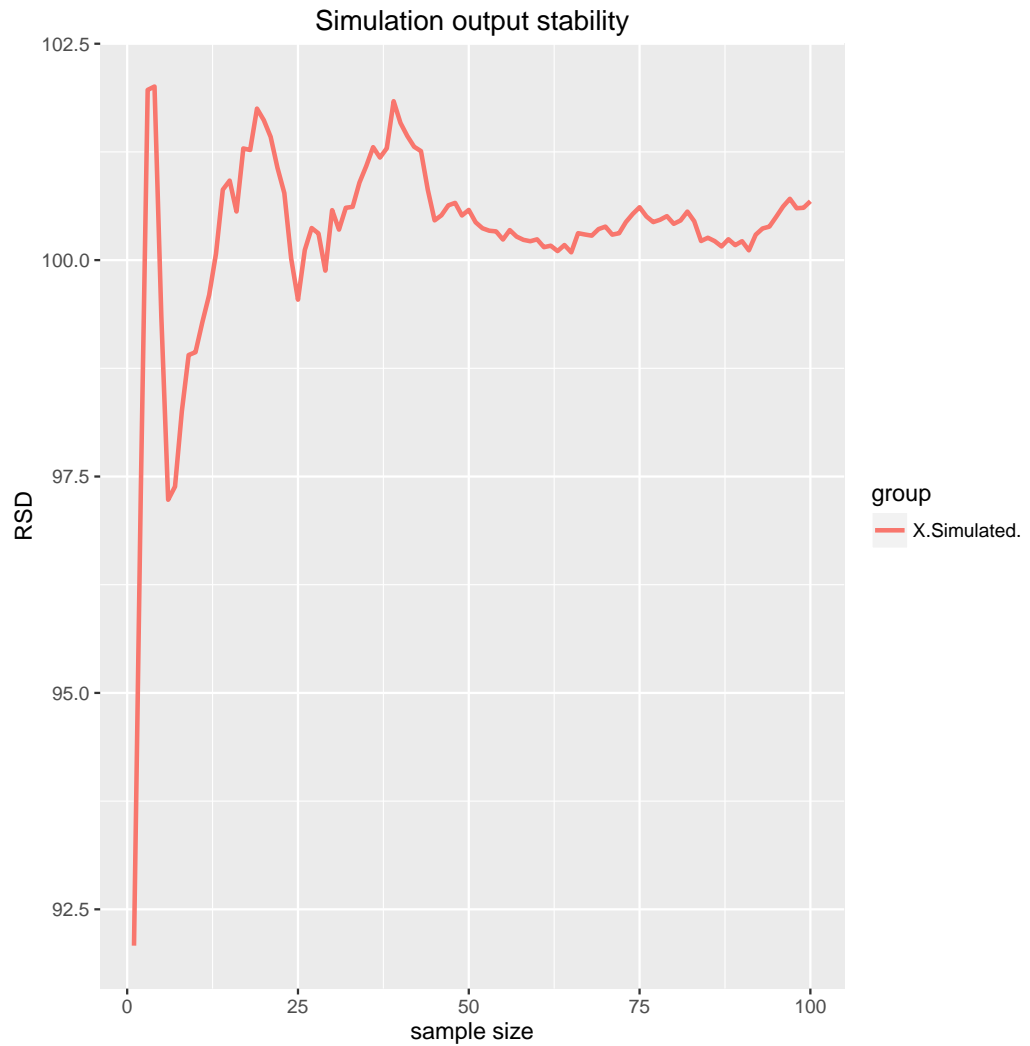
### 558 Process overview and scheduling

559 The agents are defined by the execution of a set of processes depicting the agent movement and search  
560 of food source, the consumption of food, the process incrementing the agent reserves, the reproduction  
561 and finally the death process driven by predation or starvation. The fundamental idea behind the model  
562 formulation is that both predator and prey individuals incrementing their "energy" levels by predation or  
563 by consuming the available grass respectively. Both agent types search for their food in the current patch  
564 where they are placed. The agents move a unit of space at time selecting randomly the heading.

565 The individual-based version of this model is a spatially explicit representation and have a few  
566 parameters more but is still very succinct. The list of model parameters are shown in Table 5.

<sup>4</sup>In the currently API version there is a function for accessing the charts for the *Easy.Stability* method, named *Easy.getChart()*, please refer to the user manual for the complete syntax.

<sup>5</sup>It is possible to limit it passing to the method a subset of model outputs



**Figure 7.** The stability of model output. It is possible to observe how, insofar that the number of replications of the experimental setup increases, the value of the coefficient of variation converges to a common value.

567

**Table 5.** The input parameter collection for the Repast implementation of Predator-Prey model.

Input parameter	Description
initialnumberofwolves	The initial population of predators.
initialnumberofsheep	The initial population of preys.
wolfgainfromfood	The rate of predator energy is incremented every time a prey is consumed.
wolfreproduce	The reproduction rate of predator individual.
sheepgainfromfood	The prey rate energy increment for grazing grass.
sheepreproduce	The reproduction rate of prey individual.
grassregrowthtime	The amount of time required for grass be available again once consumed by a prey.

568

569

The original formulation of Lotka-Volterra consists in a system of two differential equations with four parameters, namely the predator and the prey growth rate, the effect of predator on the prey growth and

570

571 finally the effect of prey on the predator growth as can be seen in Equation (9).

$$\begin{aligned}\frac{dx}{dt} &= c_1x - c_3xy \\ \frac{dy}{dt} &= -c_2y + c_4xy.\end{aligned}\tag{9}$$

572 There is a conceptual correspondence between the predator  $c_2$  and prey  $c_1$  growth rates with the model  
573 parameters *wolfreproduce* and *sheepreproduce* as well as with between the parameter *wolfgainfromfood*  
574 and the constant  $c_4$ .

### 575 **Model analysis**

576 The implementation code for the Morris screening exercise is shown in Figure 8 and, as has been  
577 mentioned in the previous example, we have the same sequence of steps, starting with the library loading  
578 and the selection of the random seed. Subsequently we define the objective function, which in this case is  
579 a very simple one consisting in the arithmetic average of the population sizes of sheep individuals and  
580 wolves. The next step is the selection of model input factors for the screening method and providing  
581 the range of variation for each them. Then, the **step 3** shows the call to the *Easy.Setup* function which  
582 initializes the Repast Model with the R/Repast integration code. Finally, the function *Easy.Morris*  
583 is called and the results stored in the variable *r*. The example uses five levels with ten sampling points for  
584 Morris method. The results consist in a R list holding the experiment carried out, the Morris object and a  
585 list with charts generated by the experiment<sup>6</sup>.

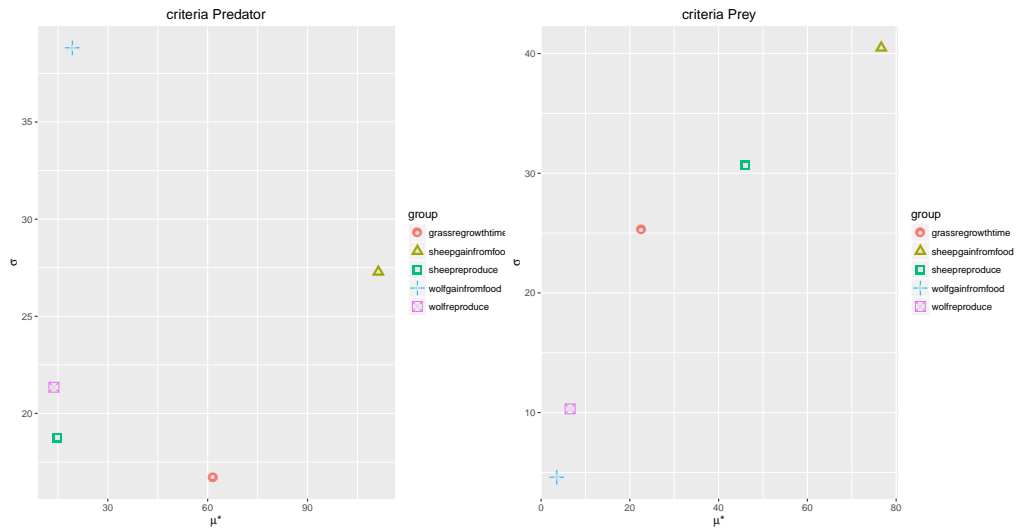
```
1 # Step 0
2 rm(list=ls())
3 library(rrepast)
4 set.seed(161803398)
5
6 # Step 1
7 objective<- function(params, results) {
8   Predator<- mean(results$Wolf.Count)
9   Prey<- mean(results$Sheep.Count)
10  cbind(Predator, Prey)
11 }
12
13 # Step 2
14 f<- AddFactor(name="wolfreproduce",min=2,max=8)
15 f<- AddFactor(factors=f, name="wolfgainfromfood",min=15,max=25)
16 f<- AddFactor(factors=f, name="sheepreproduce",min=2,max=6)
17 f<- AddFactor(factors=f, name="sheepgainfromfood",min=2,max=6)
18 f<- AddFactor(factors=f, name="grassregrowthtime",min=20,max=40)
19
20 # Step 3
21 Easy.Setup("/usr/models/PredatorPrey")
22
23 # Step 4
24 r<- Easy.Morris("/usr/models/PredatorPrey","Agent Counts",300,f,5,10,1,objective)
```

**Figure 8. The listing for Morris screening method using the *Easy.Morris* function from R/Repast.**

586 The Figure 9 presents the  $\mu^*$  vs  $\sigma$  chart for both predator and prey average population sizes. At a first  
587 glance, the most important input factor for both Predator and Prey populations is the *sheepgainfromfood*.  
588 The second most significant for the Predator output is *grassregrowthtime*. The other parameters are not  
589 very significant for the average of Predator individuals. It is also interesting to note that *wolfgainfromfood*  
590 has very high value of  $\sigma$  which could indicate that the parameter significance strongly depends on  
591 the values of other parameters. On the other hand, it could mean that the number of sampling points  
592 or replications should be increased. The Prey output presents three important parameters, which in  
593 order of importance are the *sheepgainfromfood*, the *sheepreproduce* and *grassregrowthtime*. These input  
594 parameters also have a high  $\sigma$  values which possibly indicate some non-linear effects or that the values

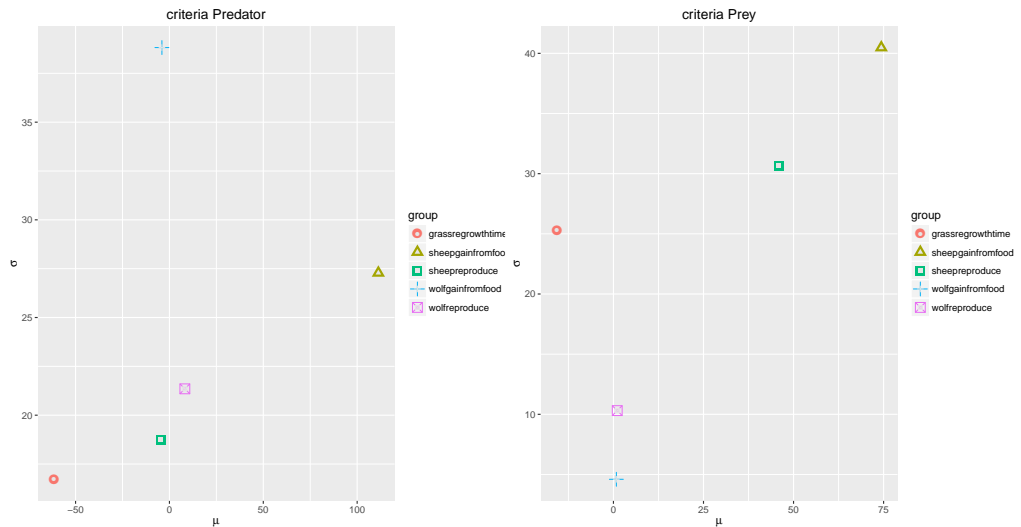
<sup>6</sup>In order to plot the charts the user should use a R code for accessing the chat list members. There are three members, namely *mustar*, *musigma* and *mumu*. In order to get the *mumu* chart for the second objective function output we must use the R call: `r$charts[2,]$mumu`.

595 of these input factors are influencing each other. These results can be explained by the dependence of  
 596 wolf population on the availability of prey. The common observed pattern in that kind of model is the  
 597 population of predators lagging in phase behind the prey population.



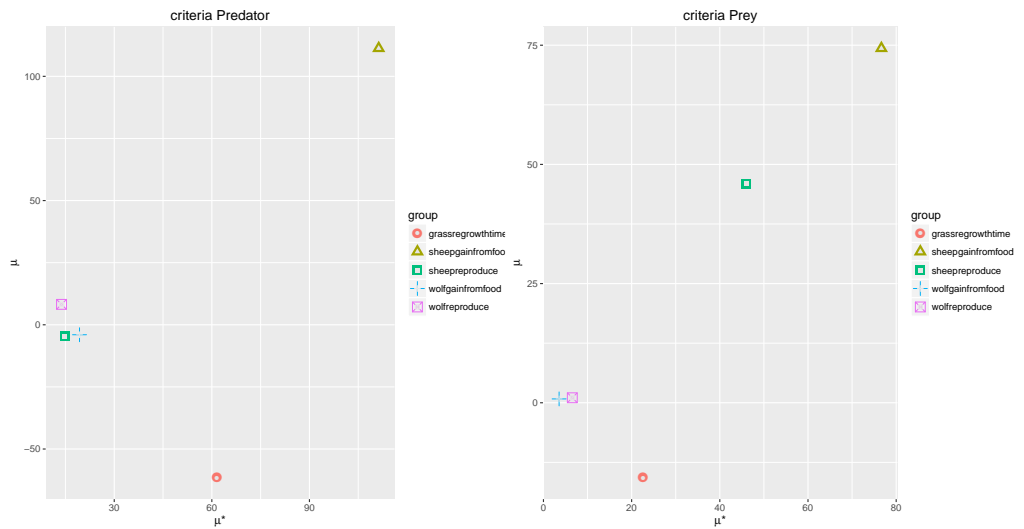
**Figure 9.** Results of Morris screening method for Predator-Prey model. The graph shows the  $\mu^*$  and  $\sigma$  sensitivity measures for Predator and Prey outputs.

598 The chart of  $\mu$  vs  $\sigma$  for model output is shown in Figure 10. It seems to provide very similar results  
 599 and the only significant difference is the contribution of *grassregrowthtime*. That input parameter was  
 600 considered important by  $\mu^*$  vs  $\sigma$  but here it has a negative value. In order to interpret this sensitivity  
 601 measure, we must recall that  $\mu^*$  takes the absolute values of elementary effects. Therefore, the elementary  
 602 effects of *grassregrowthtime* possibly has effect of opposite signs depending on the values of that input  
 603 parameter.



**Figure 10.** Results of Morris screening method for Predator-Prey model. The graph shows the  $\mu$  and  $\sigma$  sensitivity measures for Predator and Prey outputs.

604 Finally, we have the 11 showing the chart of  $\mu^*$  vs  $\mu$  where the value of both measures can be  
 605 observed together allowing the appreciation of the differences of both, which possibly indicates that the  
 606 input factors present effects with different signs which, in other words, means non-linearity in the model  
 607 behavior.



**Figure 11.** Results of Morris screening method for Predator-Prey model. The graph shows the  $\mu^*$  and  $\mu$  sensitivity measures for Predator and Prey outputs.

### 608 EXAMPLE 3: T4SS COMMON POOL

#### 609 Model description

##### 610 Purpose

611 The objective of this model is to explore the conditions where two plasmids can coexist in a population  
 612 competing for a common resource required for their horizontal transfer. The common resource is the set  
 613 of genes required for conjugation because one of the two plasmids have lost these genes.

##### 614 Entities, State variables and scales

615 The model uses two entities types, namely the agents representing the bacterial cells and a *ValueLayer*,  
 616 which is a Repast specific structure, for holding the nutrient available for the bacterial growth. The agents  
 617 interact and grow in a computational domain of  $100 \times 100 \mu m$  squared lattice with periodic boundaries  
 618 representing a total real surface of  $0.01 mm^2$ . Despite of being a lattice the bacterial cells are positioned  
 619 and allowed to move in a continuous space system. The agents are also allowed to overlap to each other.  
 620 Explicitly, the agents are defined by the five state variables: (a) heading, (b) mass, (c) division mass, (d)  
 621 plasmid *PI* infection state and (f) plasmid *PI* infection state. The current position of every bacterial cell  
 622 in the coordinate system is available implicitly through a Repast API call.

##### 623 Process overview and scheduling

624 Every bacterial cell in this model is abstracted as the execution of a series of successive processes capturing  
 625 the basic tenets of bacterial life-cycle. These processes are the nutrient uptake, the bacterial cell growth,  
 626 the division and the conjugation. The input parameters required for initializing the model are shown in  
 627 Table 6.

628 **Table 6.** The input parameter collection for the conjugative plasmid common pool model.

Input parameter	Description
doublingTime	The doubling time of plasmid free cells
p1P ( $P(\gamma_0)$ )	The probability of cell conjugate at least one time.
p1Cost	The cost imposed by the plasmid P1 including the metabolic burden required to express the conjugative apparatus.
p2Cost	The metabolic cost of plasmid P2

## 630 **Design concepts**

631

632 **Basic Principles** — The plasmid dispersion depends on an intricate balance between metabolic costs  
633 associated to horizontal and vertical dispersion strategies. The conjugative proficiency requires the  
634 expression of set of transmembrane proteins which are known collectively as Type IV Secretion Systems  
635 (T4SS) Lawley et al. (2003). The presence of conjugative plasmids and the expression of conjugative  
636 machinery is detrimental for the host cell fitness Rozkov et al. (2004) but there is no consensus on the  
637 valid ranges of metabolic costs imposed by the conjugative process. Therefore, in this model the short  
638 term dynamics of two plasmid system  $P1$  and  $P2$  is simulated. The plasmid  $P1$  is a complete conjugative  
639 plasmid containing all genes required for horizontal transfer and the plasmid  $P2$  is a cheater, which having  
640 lost its conjugative genes, depends on the T4SS system from the plasmid  $P1$ . In other words, the model  
641 is used to assess how large should be the cost difference required for the lack of conjugative apparatus  
642 become a true competitive advantage making  $P2$  dominate over  $P1$ .

643 **Emergence** — The colony growth pattern, the population distribution and the dominance of a plasmid  
644 over another on the bacterial population are global properties arising from local properties defining the  
645 agent behavior and the interaction constraints.

646 **Adaptation** — All agents adapt their growth rate, as well as the conjugation rates, according to the local  
647 availability of nutrient.

648 **Fitness** — The bacterial cells infected by any plasmid are considered to behave less efficiently than  
649 the plasmid free cells. The fitness of plasmid bearing cells are explicitly specified by the cost input  
650 parameters.

651 **Objectives** — No objectives are taken into account in this model.

652 **Prediction** — The model will provide predictions on the possible ranges of plasmid metabolic cost which  
653 can favorable to the cheaters plasmid strategy.

654 **Sensing** — The agents representing the virtual bacterial cells sense the environment to the extent that the  
655 nutrient availability controls the growth and the conjugation rates.

656 **Interaction** — Bacterial cells interact with their nearby individuals for nutrient access, cellular division,  
657 mate pair formation and plasmid transfer.

658 **Stochasticity** — Stochasticity is introduced at individual level for all cellular process sampling a normal  
659 deviate and fitting the value to corresponding process.

660 **Collectives** — No collectives are taken into account in this model.

661 **Observation** — The model provides two kind of outputs, one is numeric and contains the total number  
662 of bacterial cells which are plasmid free or are infected by the plasmids  $P1$ ,  $P2$  or both. These outputs  
663 are generated for every time step. The model also has an 2D view of colony growth updated every time tick.

664

## 665 **Model analysis**

666 The global sensitivity analysis using the Sobol variance decomposition method for the T4SS Common  
667 Pool model is shown in Figure 12. We can observe the same sequence of steps which has been previously  
668 mentioned. The objective function is defined for the average values of the model outputs named  $P1$ ,  
669  $P2$  and Both. These variable are respectively the bacterial population size infected by the  $P1$  plasmid,  
670 infected by the cheater plasmid  $P2$  and finally the number of individuals infected by both plasmids. The  
671 Sobol sensitivity indices will provide the measures of the importance of every input parameter shown in  
672 **step 2** of Figure 12 with respect to the results returned by the objective function, that is to say, the average  
673 population sizes.

674 The Figure 13 shows the first and total order indices for the model output  $P1$ . That output is the  
675 average number of bacterial cells infected just by the plasmid  $P1$ . As can be observed the most important  
676 input parameter is the bacterial cell doubling time followed by the probability  $P(\gamma_0)$ . This is an expected  
677 result as the rule for the conjugative transfer requires the bacterial cells have achieved a value rounding  
678 the 70% of cell mass at division. Other interesting aspect to note is the negative values of first order  
679 indices. Obviously the sensitivity indices should not be negative. This is consequence of a small sample  
680 size and to correct the problem we must increase it. The other important input factors for the plasmid  $P1$   
681 output are, in order of importance, the probability  $P(\gamma_0)$ , the cost of plasmid  $P2$  and the cost of plasmid  
682  $P1$ , both with similar sensitivity indices.

683 The first and total order indices for the model output showing average population size of plasmid  $P2$

```

1 # Step 0
2 rm(list=ls())
3 library(rrepast)
4 set.seed(161803398)
5
6 # Step 1
7 objective<- function(params, results) {
8   Both<- mean(results$Both)
9   P1<- mean(results$P1)
10  P2<- mean(results$P2)
11  cbind(P1, P2, Both)
12 }
13
14 # Step 2
15 f<- AddFactor(name="doublingTime",min=20,max=240)
16 f<- AddFactor(factors=f,name="p1P",min=0.1,max=0.8)
17 f<- AddFactor(factors=f,name="p1Cost",min=1,max=100)
18 f<- AddFactor(factors=f,name="p2Cost",min=1,max=100)
19
20 # Step 3
21 Easy.Setup("/usr/models/PoolT4SS")
22
23 # Step 4
24 r<- Easy.Sobol("/usr/models/PoolT4SS","output",720,f,100,20,1,objective)

```

**Figure 12.** The listing for Sobol GSA variance decomposition method using the *Easy.Sobol* function from R/Repast.

684 can be seen in Figure 14. It is possible to appreciate again, that the sensitivity indices show that the most  
685 important factor is the length of cellular cycle. The reason is simple, and can be attributed to the fact that  
686 plasmid P2 alone is only transferred vertically and depends on the plasmid P1 for horizontal transmission,  
687 being both aspects related to the cell cycle. Following in importance the doubling time we have the cost  
688 of plasmid P1, the cost of plasmid P2 and the probability  $P(\gamma_0)$ , being the sensitivity index of P1 cost,  
689 noticeably higher than the other two indices. This could be attributed probably because the plasmid P2  
690 require a significant cost difference in order to outcompete the plasmid P1 which vertical transfer.

691 Finally, in the Figure 15 we have the sensitivity indices for the output of model accounting for the  
692 average population size of bacterial cells infected by both plasmids. The importance of factors is consistent  
693 with the explanations for the previous sensitivity indices. Again, the most important model parameter is  
694 the doubling time of bacterial cell followed by the P1 and P2 cost parameters and by the probability  $P(\gamma_0)$ .

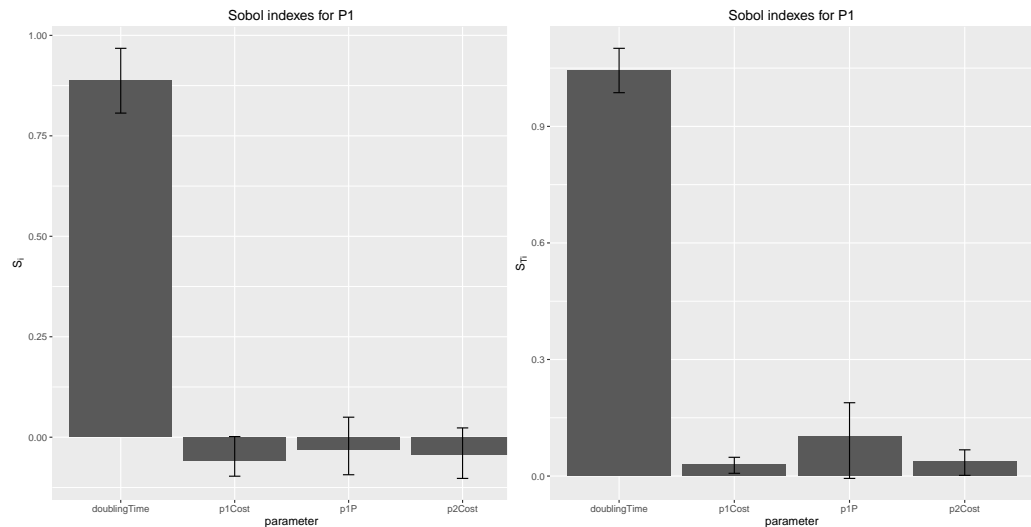
## 695 CONCLUSIONS

696 The ecological modeling is a complex subject which can be normally perceived as being simpler than it  
697 actually is. Specifically, the Individual-based models are subject to many levels of uncertainty, which  
698 means that it is hard to get completely fixed the values of model inputs, the model structure and the  
699 outputs. Normally there is no complete experimental or observational data to construct mechanistic  
700 descriptions of individual and therefore many assumptions and simplifications must be made in order to  
701 implement a model. The same is true regarding the input values, which are particularly critical in the case  
702 of the ecology of microorganism, as normally just very few input parameters are directly observed and  
703 the most of them are estimated from whole population experiments. Therefore, it is always important  
704 to bearing mind that modeling is an iterative task which must incorporate compulsorily some what-if  
705 analysis of model outputs.

706 Several methods exist for assessing the uncertainty and for estimating the relative importance of input  
707 parameters in the model output. We have provided here and overview on those methods which are based  
708 on the variance decomposition because they have a wider application scope and are specifically suitable for  
709 their use on individual-based models. These methods, although conceptually simple, are computationally  
710 intensive and can be somewhat hard to apply because the required tools are either unavailable or they do  
711 not provide an easy integration pattern. Roughly speaking, the sensitivity analysis methods require the  
712 generation of large sample of the parameter space and the model evaluation for each of them which, of  
713 course, makes the manual execution an infeasible option.

714 The in silico experimentation is becoming a vital tool for understanding complex phenomena in  
715 a way that cannot be done without modeling. The effective application of computational ecology





**Figure 13. Results of Sobol variance decomposition method for T4SS Common Pool model. The graph shows the first and total order indices sensitivity measures for bacterial population infected by plasmid P1.**

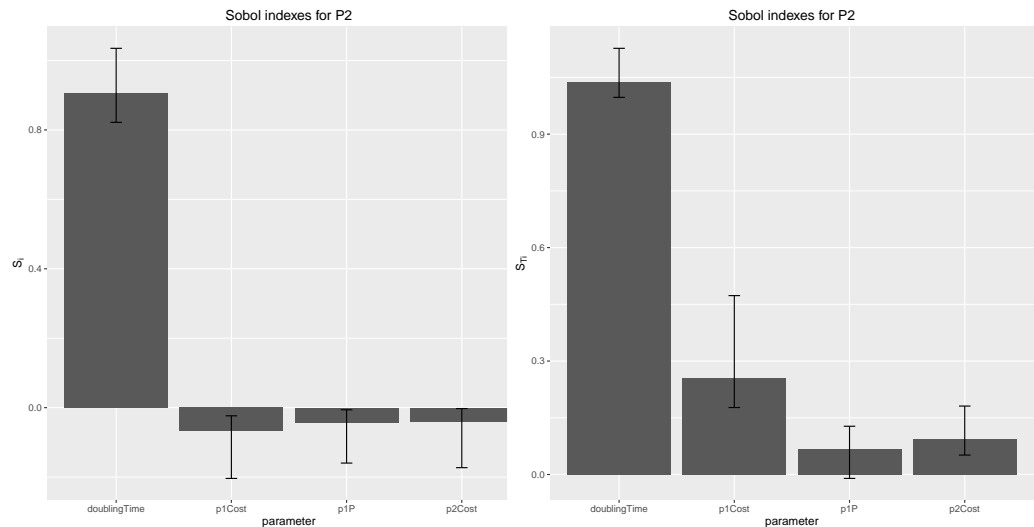
716 methods requires a high level of proficiency in many diverse domains of knowledge which sometimes  
 717 are neither feasible nor practical. Therefore, it is indispensable to have a ready to use arsenal of reusable  
 718 computational tools for modeling and analysis. In this work we have introduced the R/Repast package  
 719 and shown how it can help modelers to improve the robustness and quality of individual-based models  
 720 results by using the functionalities inside the package for analyzing systematically the model outputs. The  
 721 package can save a lot of effort for modelers by providing simple wrappers for complex methods within  
 722 a simple and consistent API. We hope that these R/Repast functionalities can facilitate enormously the  
 723 systematic analysis of Individual-based models implemented in Repast.

## 724 ACKNOWLEDGMENTS

725 This work was supported by the European FP7 - ICT - FET EU research project: 612146 (PLASWIRES  
 726 "Plasmids as Wires" project) [www.plaswires.eu](http://www.plaswires.eu) and by Spanish Government (MINECO) research  
 727 grant TIN2012-36992.

## 728 REFERENCES

- 729 Andres, T. and Hajas, W. (1993). Using iterated fractional factorial design to screen parameters in  
 730 sensitivity analysis of a probabilistic risk assessment model.
- 731 Arutyunov, D. and Frost, L. S. (2013). F conjugation: Back to the beginning. *Plasmid*, 70(1):18–32.
- 732 Beck, J. V. and Arnold, K. J. (1977). *Parameter estimation in engineering and science*. Wiley series in  
 733 probability and mathematical statistics. Wiley, New York.
- 734 Bercé, L. (2002). Techniques of spatially explicit individual-based models: construction, simulation, and  
 735 mean-field analysis. *Ecological Modelling*, 150(1-2):55–81.
- 736 Bergstrom, C. T., Lipsitch, M., and Levin, B. R. (2000). Natural Selection, Infectious Transfer and the  
 737 Existence Conditions for Bacterial Plasmids. *Genetics*, 155(4):1505–1519.
- 738 Bettonvil, B. and Kleijnen, J. P. C. (1996). Searching for important factors in simulation models with many  
 739 factors: Sequential bifurcation. *EUROPEAN JOURNAL OF OPERATIONAL RESEARCH ELSEVIER*  
 740 *European Journal of Operational Research*, 96:180–194.
- 741 Bonabeau, E. (2002). Agent-based modeling: Methods and techniques for simulating human systems.  
 742 *Proceedings of the National Academy of Sciences*, 99(Supplement 3):7280–7287.
- 743 Box, G. E. P. and Draper, N. R. (1987). *Empirical Model-Building and Response Surfaces (Wiley Series*  
 744 *in Probability and Statistics)*. Wiley, 1 edition.



**Figure 14. Results of Sobol variance decomposition method for T4SS Common Pool model. The graph shows the first and total order indices sensitivity measures for bacterial population infected by plasmid P2.**

745 Campolongo, F., Cariboni, J., and Saltelli, A. (2007). An effective screening design for sensitivity analysis  
746 of large models.

747 Chen, I., Christie, P. J., and Dubnau, D. (2005). The ins and outs of DNA transfer in bacteria. *Science*  
748 (*New York, N.Y.*), 310(5753):1456–1460.

749 Crawley, M. J. (2007). *The R Book*. Wiley, 1 edition.

750 Dieckmann, U., Law, R., and Metz, J. A. J. (2000). *The Geometry of Ecological Interactions*, volume 1.

751 Emrich, Š., Suslov, S., and Judex, F. (2007). Fully agent based modellings of epidemic spread using  
752 Anylogic. *Proceeding EUROSIM 2007, Ljubljana, Slovenia, 2007:1–7*.

753 Evans, M. R., Grimm, V., Johst, K., Knuuttila, T., et al. (2013). Do simple models lead to generality in  
754 ecology? *Trends in ecology & evolution*, 28(10):578–83.

755 Ferrer, J., Prats, C., and López, D. (2008). Individual-based modelling: an essential tool for microbiology.  
756 *Journal of biological physics*, 34(1-2):19–37.

757 Grimm, V., Berger, U., Bastiansen, F., Eliassen, S., Ginot, V., Giske, J., Goss-Custard, J., Grand, T., Heinz,  
758 S. K., Huse, G., Huth, A., Jepsen, J. U., Jørgensen, C., Mooij, W. M., Müller, B., Pe'er, G., Piou, C.,  
759 Railsback, S. F., Robbins, A. M., Robbins, M. M., Rossmanith, E., Rüger, N., Strand, E., Souissi, S.,  
760 Stillman, R. A., Vabø, R., Visser, U., and DeAngelis, D. L. (2006). A standard protocol for describing  
761 individual-based and agent-based models. *Ecological Modelling*, 198(1-2):115–126.

762 Grimm, V., Berger, U., DeAngelis, D. L., Polhill, J. G., Giske, J., and Railsback, S. F. (2010). The ODD  
763 protocol: A review and first update. *Ecological Modelling*, 221(23):2760–2768.

764 Grimm, V. and Railsback, S. F. (2005). *Individual-based Modeling and Ecology: (Princeton Series in*  
765 *Theoretical and Computational Biology)*. Princeton University Press, Princeton.

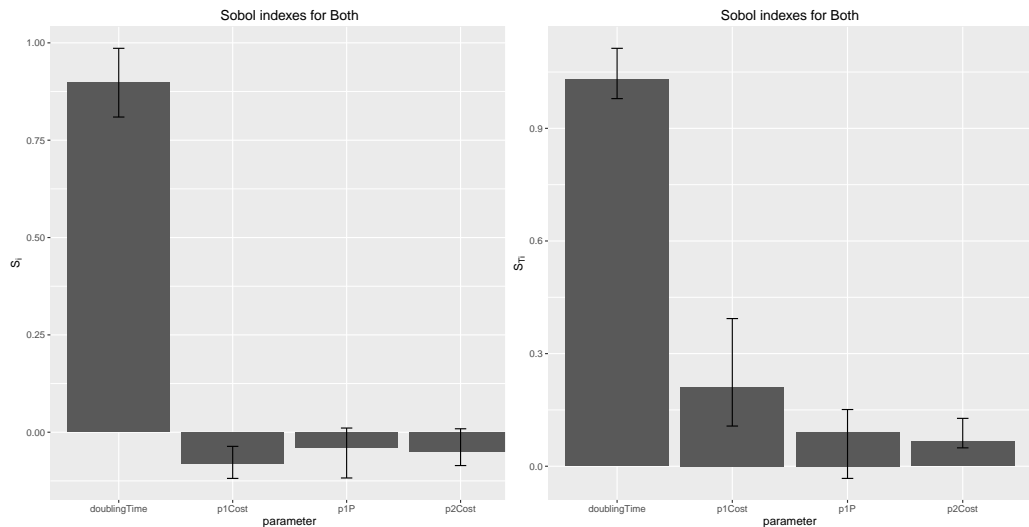
766 Gutfraind, A., Boodram, B., Prachand, N., Hailegiorgis, A., Dahari, H., and Major, M. E. (2015). Agent-  
767 Based Model Forecasts Aging of the Population of People Who Inject Drugs in Metropolitan Chicago  
768 and Changing Prevalence of Hepatitis C Infections. *PLOS ONE*, 10(9):e0137993.

769 Happe, K., Kellermann, K., and Balmann, A. (2006). Agent-based analysis of agricultural policies: An  
770 illustration of the agricultural policy simulator AgriPolis, its adaptation and behavior. *Ecology and*  
771 *Society*, 11(1).

772 Hellweger, F. L. and Bucci, V. (2009). A bunch of tiny individuals — Individual-based modeling for  
773 microbes. *Ecological Modelling*, 220(1):8–22.

774 Helly, J., Case, T., Davis, F., Levin, S., and Michener, W. (1995). *The State of Computational Ecology*.  
775 *San Diego, CA: San Diego Super Computer Center*.

776 Herman, J. D., Kollat, J. B., Reed, P. M., and Wagener, T. (2013). Technical Note: Method of Morris  
777 effectively reduces the computational demands of global sensitivity analysis for distributed watershed



**Figure 15. Results of Sobol variance decomposition method for T4SS Common Pool model. The graph shows the first and total order indices sensitivity measures for bacterial population infected by both plasmids P1 and P2.**

778 models. *Hydrology and Earth System Sciences*, 17(7):2893–2903.

779 Hicks, C. R. (1993). *Fundamental Concepts in the Design of Experiments*. Oxford University Press, USA,

780 4 edition.

781 Law, A. M. (2005). How to build valid and credible simulation models. In Kuhl, M. E., Steiger, N. M.,

782 Armstrong, F. B., and Joines, J. A., editors, *Proceedings of the 2005 Winter Simulation Conference*,

783 pages 24–32.

784 Lawley, T. D., Klimke, W. A., Gubbins, M. J., and Frost, L. S. (2003). F factor conjugation is a true type

785 IV secretion system. *FEMS microbiology letters*, 224(1):1–15.

786 Little, T. M. T. M. and Hills, F. J. (1978). *Agricultural experimentation : design and analysis*. Wiley.

787 Lorscheid, I., Heine, B.-O., and Meyer, M. (2012). Opening the 'black box' of simulations: increased

788 transparency and effective communication through the systematic design of experiments. *Computational*

789 *and Mathematical Organization Theory*, 18(1):22–62.

790 Luke, S., Cioffi-Revilla, C., Panait, L., Sullivan, K., and Balan, G. C. (2005). MASON: A Multiagent

791 Simulation Environment. *Simulation*, 81(7):517–527.

792 Morris, M. D. (1991). Factorial Sampling Plans for Preliminary Computational Experiments. *Technome-*

793 *trics*, 33(2):161–174.

794 Myers, J. L. and Well, A. D. (1995). *Research Design & Statistical Analysis*. Routledge, 1 edition.

795 North, M., Collier, N., and Ozik, J. (2013a). Complex adaptive systems modeling with repast simphony.

796 *Complex adaptive ...*, pages 1–26.

797 North, M., Collier, N., Ozik, J., Tatara, E., Macal, C., Bragen, M., and Sydelko, P. (2013b). Complex

798 adaptive systems modeling with Repast Simphony. *Complex Adaptive Systems Modeling*, 1(1):1–26.

799 North, M. J., Collier, N. T., Ozik, J., Tatara, E. R., Macal, C. M., Bragen, M., and Sydelko, P. (2013c).

800 Complex adaptive systems modeling with Repast Simphony. *Complex Adaptive Systems Modeling*,

801 1(1):3.

802 North, M. M. J., Collier, N. T. N., Ozik, J., Tatara, E. R., Macal, C. M., Bragen, M., and Sydelko, P.

803 (2013d). Complex adaptive systems modeling with Repast Simphony. *Complex Adaptive Systems*

804 *Modeling*, 1(1):1–26.

805 Pascual, M. (2005). Computational Ecology: From the Complex to the Simple and Back. *PLoS*

806 *Computational Biology*, 1(2):e18.

807 Petrovskii, S., Petrovskaya, N., Hughes, J. D., et al. (2012). Computational ecology as an emerging

808 science. *Interface focus*, 2(2):241–54.

809 Pianosi, F., Beven, K., Freer, J., Hall, J. W., Rougier, J., Stephenson, D. B., and Wagener, T. (2016). Sen-

810 sitivity analysis of environmental models: A systematic review with practical workflow. *Environmental*

811 *Modelling & Software*, 79:214–232.

812 Prestes García, A. and Rodríguez-Patón, A. (2015). A Preliminary Assessment of Three Strategies for the  
813 Agent-Based Modeling of Bacterial Conjugation. In Overbeek, R., Rocha, M. P., Fdez-Riverola, F.,  
814 and De Paz, J. F., editors, *9th International Conference on Practical Applications of Computational  
815 Biology and Bioinformatics*, volume 375 of *Advances in Intelligent Systems and Computing*, pages 1–9.  
816 Springer International Publishing.

817 Prestes García, A. and Rodríguez-Patón, A. (2015). BactoSim — An Individual-Based Simulation  
818 Environment for Bacterial Conjugation. pages 275–279. Springer International Publishing.

819 Pujol, G., Iooss, B., with contributions from Sebastien Da Veiga, A. J., Fruth, J., Gilquin, L., Guillaume,  
820 J., Gratiot, L. L., Lemaitre, P., Ramos, B., and Touati, T. (2015). *sensitivity: Sensitivity Analysis*. R  
821 package version 1.11.1.

822 R Core Team (2015). *R: A Language and Environment for Statistical Computing*. R Foundation for  
823 Statistical Computing, Vienna, Austria.

824 Rozkov, A., Avignone-Rossa, C. A., Ertl, P. F., Jones, P., O’Kennedy, R. D., Smith, J. J., Dale, J. W., and  
825 Bushell, M. E. (2004). Characterization of the metabolic burden on *Escherichia coli* DH1 cells imposed  
826 by the presence of a plasmid containing a gene therapy sequence. *Biotechnology and bioengineering*,  
827 88(7):909–915.

828 Saltelli, A. (2008). Global Sensitivity Analysis: The Primer. *International Statistical Review*, page 452.

829 Saltelli, A., Tarantola, S., Campolongo, F., and Ratto, M. (2004). *Sensitivity Analysis in Practice: A  
830 Guide to Assessing Scientific Models*. Wiley, 1 edition.

831 Seoane, J., Yankelevich, T., Dechesne, A., Merkey, B., Sternberg, C., and Smets, B. F. (2011). An  
832 individual-based approach to explain plasmid invasion in bacterial populations. *FEMS microbiology  
833 ecology*, 75(1):17–27.

834 Slater, F. R., Bailey, M. J., Tett, A. J., and Turner, S. L. (2008). Progress towards understanding the fate of  
835 plasmids in bacterial communities. *FEMS Microbiology Ecology*, 66(1):3–13.

836 Smith, J. M. (1974). *Models in ecology*. Cambridge University Press, revised ed edition.

837 Tack, I. L., Logist, F., Noriega Fernández, E., and Van Impe, J. F. (2015). An individual-based modeling  
838 approach to simulate the effects of cellular nutrient competition on *Escherichia coli* K-12 MG1655  
839 colony behavior and interactions in aerobic structured food systems. *Food Microbiology*, 45:179–188.

840 Thiele, J. C., Kurth, W., and Grimm, V. (2014). Facilitating Parameter Estimation and Sensitivity Analysis  
841 of Agent-Based Models: A Cookbook Using NetLogo and ‘R’. *Journal of Artificial Societies and  
842 Social Simulation*, 17(3).

843 Tisue, S. and Wilensky, U. (2004). Netlogo: A simple environment for modeling complexity. ...  
844 *Conference on Complex Systems*, pages 1–10.

845 Urbanek, S. (2016). *rJava: Low-Level R to Java Interface*. R package version 0.9-8.

846 van Houwelingen, H. C., Boshuizen, H. C., and Capannesi, M. (2011). Sensitivity analysis of state-  
847 transition models: How to deal with a large number of inputs. *Computers in Biology and Medicine*,  
848 41(9):838–842.

849 Watkins, A., Noble, J., Foster, R., Harmsen, B., and Doncaster, C. (2015). A spatially explicit agent-  
850 based model of the interactions between jaguar populations and their habitats. *Ecological Modelling*,  
851 306:268–277.

852 Xu, C. and Gertner, G. (2011). Understanding and comparisons of different sampling approaches for the  
853 Fourier Amplitudes Sensitivity Test (FAST). *Computational statistics & data analysis*, 55(1):184–198.

854 Zeigler, B. P., Praehofer, H., and Kim, T. G. (2000). *Theory of Modeling and Simulation, Second Edition*.  
855 Academic Press, 2 edition.

856 Zhang, Y. and Rundell, A. (2006). Comparative study of parameter sensitivity analyses of the TCR-  
857 activated Erk-MAPK signalling pathway. *Systems biology*, 153(4):201–11.

858 Zhong, X., Droesch, J., Fox, R., Top, E. M., and Krone, S. M. (2012). On the meaning and estimation  
859 of plasmid transfer rates for surface-associated and well-mixed bacterial populations. *Journal of  
860 Theoretical Biology*, 294:144–152.