

A generic multi-level stochastic modelling framework in computational epidemiology

Sébastien Picault^{1,2*}, Yu-Lin Huang¹, Vianney Sicard¹, Thierry Hoch¹, Elisabeta Vergu³, François Beaudeau¹, Pauline Ezanno¹

¹BIOEPAR, INRA, Oniris, CS40706, 44307 Nantes, France

²Univ. Lille, CNRS, Centrale Lille, UMR 9189 - CRISTAL, Lille, France

³MaIAGE, INRA, Université Paris-Saclay, Jouy-en-Josas, France

*Corresponding author's Email: sebastien.picault@oniris-nantes.fr

Abstract

There is currently an overwhelming increased interest in predictive biology and computational modelling. The development of reliable, reproducible and revisable simulation models in computational life sciences is often pointed out as a challenging issue. Population dynamics, including epidemiology, has not yet developed a language to formalize complex models in a univocal and automatable way, hence hindering the capability to implement in short time reliable, revisable and expert-friendly models intended for realistic mechanistic simulations. In epidemiology specifically, models aim not only at understanding pathogen spread but also at assessing control measures at several scales. To achieve this goal efficiently, best software practices should be supported by Artificial Intelligence methods to handle experts' knowledge. The framework EMULSION presented here intends to both tackle multiple modelling paradigms in epidemiology and facilitate the automation of model design. We therefore built both a domain-specific language (DSL) for the modular description of complex epidemiological models, and a generic simulation engine designed to embed existing modelling paradigms within a homogeneous architecture based on adaptive software agents. The diversity of concerns (biology, economics, human activities) involved in real pathosystems requires an explicit, comprehensive and intelligible way to describe epidemiological models, to involve experts without computer science skills throughout the modelling, simulation and output analysis steps. This approach was applied to compare hypotheses in modelling a zoonosis (Q fever), to study its transmission dynamics within and between cattle herds at a regional scale, and to assess the contribution of transmission pathways. Separating model description from the simulation engine allowed epidemiologists to be involved in assumption revision, while guaranteeing very few code modifications. We assessed the added value of EMULSION by applying the DSL and the simulation engine to a concrete disease. Future extensions of EMULSION towards a broader range of epidemiological concerns will reduce significantly the time required to design and assess models and control measures against endemic and epidemic diseases. Ultimately, we believe this effort is a major lever to increase scientists' preparedness to face emerging threats for public health and provide rapid, reliable, and reasoned assessments of control measures.

Keywords

Computational Biology, Artificial Intelligence, Infectious Diseases, Epidemiological Modelling, Multi-Scale Models, Agent-Based Simulation, Knowledge and Software Engineering, Domain-Specific Language, Expert Involvement

1 **Introduction**

2 **Balancing development time, reliability and intelligibility in computational models**

3 Computational modelling is essential to better explore complex systems. In particular,
4 agricultural production systems present highly coupled biological, farming, environmental, and
5 economic processes, involving a diversity of interacting entities, from individual scale up to whole
6 territories. Their analytical investigation is strongly limited by the interplay between all processes and
7 scales, leading to high dimensional and highly nonlinear systems, but also by the boundaries of
8 knowledge concerning the exact interactions between actors of such systems. Mechanistic simulation
9 models can assess the relevance of assumptions by comparing model outputs to field data, provide
10 predictions on systems evolution under real or counterfactual scenarios, and help identify levers to
11 control those systems. However, it is crucial that alternative hypotheses and practicable actions be tested
12 in short time, in strong interaction with experts, to quickly identify assumptions providing the most
13 significant insights, or actions driving the system to a desired state. Such "sieving" of hypotheses also
14 promotes parsimonious models, highlighting key elements, hence allowing for deeper understanding
15 and easier comparisons.

16 However, developing simulation codes directly from models requires strong skills in computer
17 programming. Any change in hypotheses, scenarios, model structure or even just parameters is
18 excessively time-consuming to foster incremental design of models and expert involvement. Also,
19 reliability and reproducibility issues of ad-hoc simulation programs threaten conclusions drawn from
20 computer experiments. To avoid misinterpretations coming from programming biases [1], several good
21 practices in software development were proposed [2] (e.g. precise code documentation, systematic
22 testing, versioning, etc.), but erroneous programs can also reach such standards [3].

23 Models are intended to change with biological knowledge and research questions. Assessing
24 their relevance rather than simulation code quality requires to allow experts scrutinizing their very
25 components (parameters, functions, modelling paradigms, contact structures, etc.), instead of their
26 implementation within a specific programming language. It is then the responsibility of computer
27 scientists to provide an automated, reliable and rapid translation into code. Our approach is in line with
28 this mindset, by coupling a modular simulation architecture with a Domain-Specific Language (DSL),

29 which gives experts the ability to understand and design the multiple components of an epidemiological
30 model without programming.

31 **The diversity of modelling and related issues in epidemiology**

32 Epidemiology is an epitome field for addressing such issues. Since Kermack and McKendrick's
33 seminal works [4], the complexity of models increased to allow for realistic decision support at several
34 scales [5, 6], incorporating a broad range of concerns: infectious processes, demography, environmental
35 conditions, underlying contact structure provided by transportation or trade, etc. Models became hard
36 to design and harder to implement in a reliable way, because life scientists are not expected to master
37 programming skills [7]. The diversity of modelling paradigms, as presented below, from rather formal
38 and analytical, to rule-based descriptions of processes involved in the system, also reduces the ability to
39 revise or compare models in response to evolving scientific knowledge and purposes. This often results
40 in heterogeneous, ad-hoc simulation programs which cannot be compared, enhanced, even used, without
41 diving into the code. However, responsiveness in modelling and in scenario assessment is a stake to
42 provide relevant control measures against outbreaks, especially in the case of an animal health crisis.

43 Compartment-based models (CBM) [8] describe disease dynamics by state variables (amount or
44 proportion of individuals in each health state). CBM can also represent demographic dynamics with
45 input and output rates, and incorporate additional concerns (e.g. age structure, species, or environment-
46 borne contamination) by splitting compartments. CBM assume that individuals differ only by a few
47 discrete variables which determine their compartment. To assess targeted control measures, the
48 multiplication of compartments and transitions required to account for finer-grained features can make
49 the model very like individual-based models (IBM) [9]. These latter keep individual diversity explicit
50 [10, 11] and represent them with their behaviors, environment, possible goals (e.g. [12–16]). This
51 comprehensive understanding of causal mechanisms occurring in biological systems allows to compare
52 individual trajectories and measure the impact of fine-grained actions [17]. The capability to increase
53 indefinitely detail level as needed is counterbalanced by a high computational cost and by a difficulty
54 to calibrate parameters (even with parsimonious models), two major drawbacks of IBM, since the
55 scientific soundness of simulation outcomes relies upon repetitions and sensitivity analysis. Their use
56 on very large scales (e.g. millions of agents) is a challenge, even with massively parallel platforms,
57 strong software optimizations and oversimplified epidemiological assumptions [13, 18].

58 Metapopulations approaches [19, 20] have been applied in epidemiology for handling region-wide
59 models at a moderate computational cost. Populations are modelled in interaction through a contact
60 structure [21] representing neighborhood relations, transportation or trade exchanges [22, 23], or vector-
61 or airborne transport processes. Yet, approximations in sub-populations dynamics may result in
62 overestimating infections [24, 25] compared to equivalent IBM.

63 Most paradigms share the flow diagram formalism, with nodes denoting health states (possibly
64 combined with other concerns), and transitions labeled with rates. In continuous, deterministic
65 approaches, they are equivalent to an Ordinary Differential Equation (ODE) system, while in stochastic
66 models, rates can be used, after conversion into probabilities, either in discrete event methods (e.g. the
67 Gillespie algorithm [26]), or in multinomial sampling in discrete time approaches [27]. The main
68 drawback of flow diagrams is that several concerns (infection, demography, herd management...) often
69 are mixed together in a single representation, reducing the readability of the model, while other features
70 (parameters, processes, data...) are not systematically explicitly depicted (e.g. the exponential
71 distribution of state durations, or pathogen shedding during infectious states). Then, when writing actual
72 simulation code, several implicit assumptions or actions are just added on the fly, which hinders early
73 model comparison and often leads to biases in the late stage.

74 **Related computer science solutions and specificities of our approach**

75 Epidemiology does not provide any systematic methodology for designing, implementing or
76 assessing the diversity of its models yet. Other life sciences, which have long faced major computational
77 problems, have adopted powerful formalisms to express their models in a quite explicit and
78 comprehensive way and automatize their development. For instance, in molecular biology, the Systems
79 Biology Graphical Notation (SBGN) [28] offers a visual syntax for describing reactions, compounds or
80 feedback loops. In epidemiology, such attempts are still at their early stage. Formalisms inspired from
81 multi-scale processes in physics [29], or proposing a strong complexification of flow diagrams [30], are
82 not likely to facilitate the appropriation of models by epidemiologists. Conversely, the ODD protocol
83 ("Overview, Design concepts, Details", [31]) aims to obtain comprehensive knowledge from
84 disciplinary experts within a textual template; however, feedbacks on its actual use for designing models
85 emphasize ambiguities and "the lack of real specifications" [32].

86 Most simulation programs developed for implementing epidemiological models are hand-
87 written ad-hoc tools dedicated to a single pathogen in an applicative context, to evaluate a specific set
88 of control measures. Reliability of such codes inherently depends on the programming skills of their
89 authors, and prove difficult to use and maintain in the long-term. Especially, instead of high-level
90 programming languages (Scilab, R, Python...), performance considerations lead to using low-level ones
91 (C++), yet harder to master, debug and maintain. Besides, even in object-oriented development,
92 abstraction is rarely used (excepted e.g. in [33] on vector-borne disease mechanisms).

93 However, classical general-purpose simulation platforms tend progressively to be used, for
94 instance GAMA [14], NetLogo [16] or Repast [34]. They provide indeed reliable, ready-made tools for
95 calculations or data integration, leaving more time to focus on modelling itself. Yet, they are not
96 specifically designed for epidemiology and still require a significant time spent on software
97 development. Simulation libraries and platforms dedicated to epidemiological issues are rising, e.g.
98 SimInf [35], a R library for data-driven CBM; MicroSim [36], an agent-based platform for several
99 diseases; or GLEaMviz [37], a metapopulation-oriented platform. To our knowledge, the most advanced
100 approach from a software engineering viewpoint is Broadwick [38], a framework for CBM and IBM
101 with interaction networks, which still requires writing large portions of code to derive specific classes
102 and carry out simulations on practical cases. Another interesting approach, though dedicated to a specific
103 study, relies upon geographical levels and a separation of activities to define efficient aggregations of
104 individuals [39]. KENDRICK [40] defines a DSL fostering a strong separation of concerns (infectious
105 dynamics, spatial distribution, species), and generates C/C++ code to run simulations efficiently. But, it
106 only targets CBM with theoretical (SIR-type) models.

107 The framework EMULSION (for “Epidemiological Multi-Level Simulation”) we propose is to
108 our knowledge the only contribution that combines the capability of integrating several modelling
109 paradigms and several scales with dynamic aggregation levels (through a multi-level multi-agent
110 architecture to wrap them), and the explicit description of models (through a DSL dedicated to
111 epidemiological issues) [41] (Fig 1). Our objective was indeed to 1) define a formalism making models
112 as accurate as possible, so that a comprehensive description can be shared amongst experts, and its
113 implementation automatized; and 2) encompass existing modelling paradigms within a common
114 interface, to make them interchangeable, or even combine them as proposed by [42].

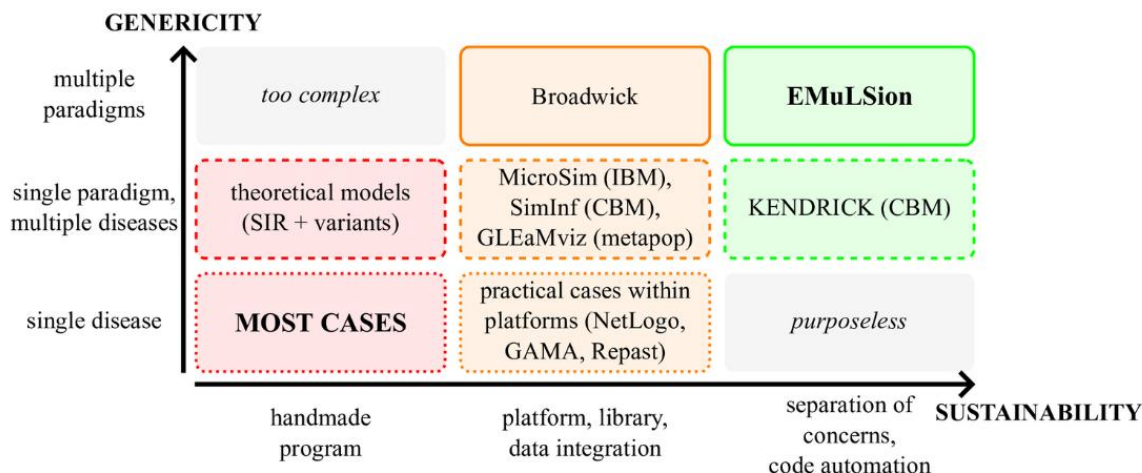


Fig 1. A taxonomy of modelling approaches in computational epidemiology. The vertical axis is based on the scope (from a single disease to multiple modelling paradigms); the horizontal axis represents the level of computer science complexity involved (from ad-hoc monolithic programs to a full separation between explicit knowledge and simulation code). In EMULSION, the coupling of a modular simulation architecture with a DSL is beneficial on both levels.

115 **Methods**

116 **Coupling a Domain-Specific Language with a generic simulation engine**

117 Designing a model using EMULSION involves three interdependent elements: 1) an explicit,
 118 modular and readable description of the model written using EMULSION’s DSL — this step is intended
 119 to be accessible to non-computer scientists experts ; 2) the use of the generic simulation engine written
 120 by computer scientists, to capitalize, in an extensible, modular and reliable way, recurrent treatments
 121 and calculations that can be found in most epidemiological models (e.g. computation of states evolution
 122 over time, connection to data, etc.) — this engine is aimed at building and running the appropriate
 123 simulation architecture based on the model specifications in the DSL; 3) small code add-ons which may
 124 be necessary to add features (calculations, actions, processes) either specific to each model or not yet
 125 incorporated into the generic engine (Fig 2). The complexity of designing and implementing a model
 126 is thus broken down into several simpler concerns, without unnecessary code writing. Besides, models
 127 described through EMULSION’s DSL are univocal in the sense that they have to make most
 128 assumptions explicit, and refer to the simulation methods implemented within the generic engine.

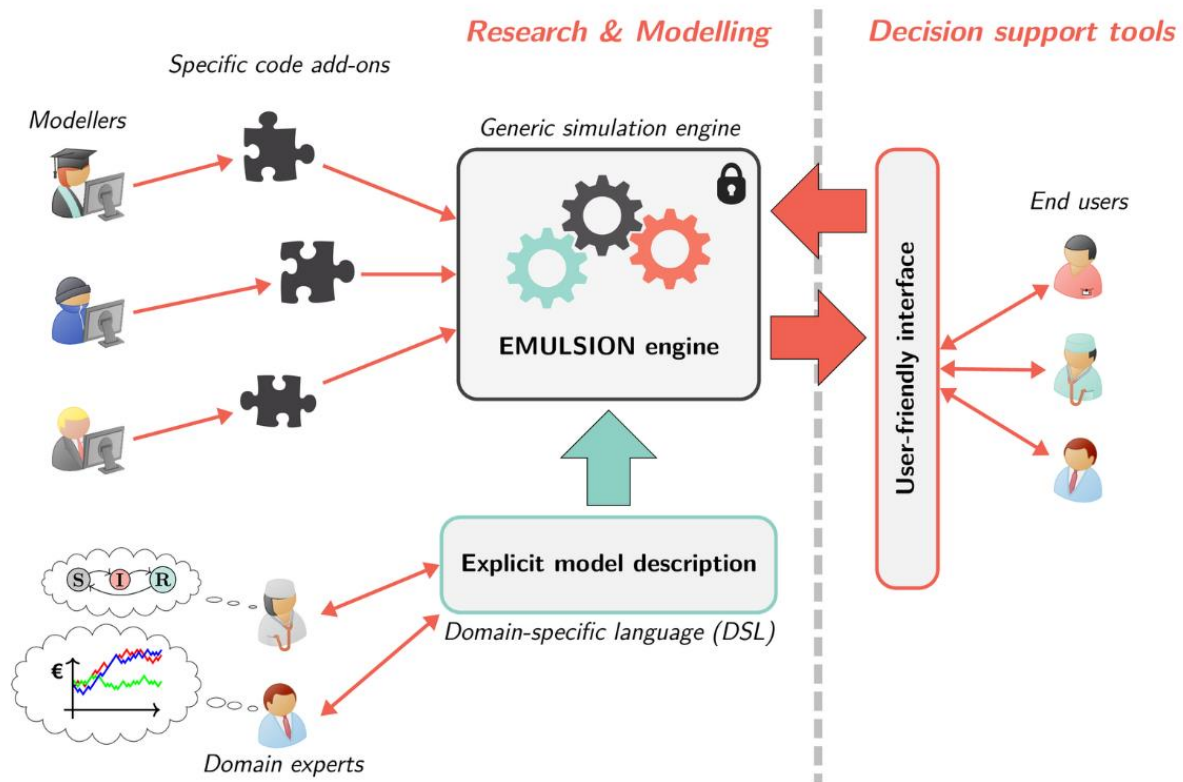


Fig 2. Approach enforced in EMULSION. A generic simulation engine is coupled to a domain-specific modelling language, fostering continuous experts' involvement and user-friendly interactions. Experts' knowledge is kept explicit, understandable, and revisable. A few specific code add-ons can be written to complement the simulation engine.

129 **Knowledge engineering: a paradigm-independent representation of processes**

130 Epidemiological models mainly rely on the description of infectious processes. As a balanced
131 formalism, we propose to extend flow diagrams to represent state evolution through Finite State
132 Machines [43], widespread used in computer science. Features that were implicit in epidemiological
133 design can be described explicitly in nodes (states) and edges (transitions) of state machine diagrams
134 (Fig 3), enhancing the intelligibility of models. States can be endowed with 1) a duration distribution,
135 specifying how long an individual is expected to stay in the current state, and 2) actions performed by
136 individuals when entering, being in, or leaving the state. Transitions are labeled with either a rate, a
137 probability or an absolute amount; they can also specify: 1) calendar conditions to indicate time periods
138 when transitions are available; 2) escape conditions allowing to free from state duration constraints; 3)
139 individual conditions to filter who is allowed; 4) actions performed by individuals crossing the transition
140 (after leaving their current state and before entering their new one).

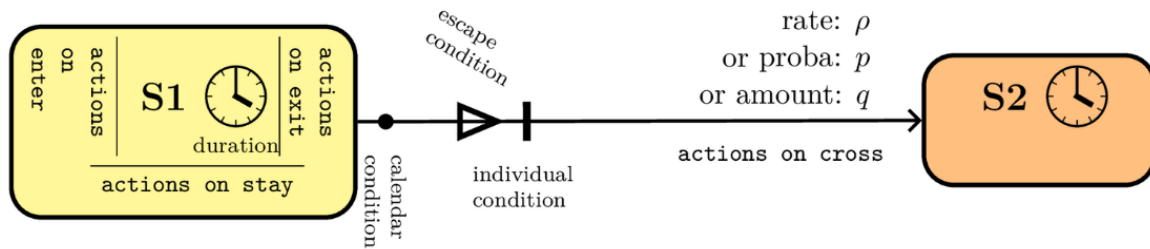


Fig 3. Structure of a transition between two states in state machines. States can be given a duration and actions when entering, staying in, or leaving the state. Transitions feature a rate, or probability, or amount, and can be associated with actions performed on crossing, time-dependent ("calendar") conditions, or individual conditions restricting the capability to cross the transition, and escape conditions allowing individuals to leave their state before the nominal duration.

141 A classical flow diagram is essentially a conceptual sketch of the model, requiring further
 142 programming to control transitions between states; conversely, the state machine diagram is informative
 143 enough to be processed directly by the generic simulation engine without further code writing. For
 144 instance, flow diagrams generally assume an exponential distribution of durations in health states; but,
 145 sometimes other distributions are required, e.g. a constant incubation duration. While a classical
 146 approach would require to rewrite the simulation code to switch from exponential to constant duration,
 147 with EMULSION the specification of a constant duration in the node of the state machine is
 148 automatically handled with the correct computation by the generic engine. Similar changes can be made
 149 or revised on the fly, since they require no more than adding or removing a few lines in the model
 150 description.

151 Each state machine is aimed at describing a single process (infection, demography, ...). Thus,
 152 instead of mixing concerns within a single, complex diagram, each process can be expressed, assessed,
 153 revised independently from the others and in a simple representation. Possible interactions between
 154 processes can be expressed through actions: for instance, actions performed in a ‘Treated’ state of a
 155 treatment process can induce changes in the infectious process.

156 While flow diagrams were population-oriented (i.e. describing the evolution of group size), state
 157 machines are individual-oriented: they specify individual behaviors, allowing to focus on fine-grained
 158 individual features. The subsequent issue, consisting in aggregating individuals to the relevant detail
 159 level without excessive computational cost, is addressed by the agent-based simulation architecture
 160 described a few lines below.

161 **A language for epidemiological model description**

162 Model assessment, from the first assumptions to the interpretation of simulation results, is a long
163 process. To keep the model explicit, understandable, and revisable throughout, it must be accessible
164 under a readable form, rather than buried into the simulation code. Thus, we recommend gathering all
165 model components (parameters, distributions, functions, time management, state machines, levels,
166 processes occurring on each scale, actions, etc.) within a structured text file. We defined a Domain-
167 Specific Language (DSL) [44] matching the needs of epidemiology, to allow experts to structure models
168 through key-value pairs. Model description files are intended to be comprehensive documents, thus
169 force modelers to leave comments and sources for each item: the same file can then be used to produce
170 figures, parameter tables, or technical documentation. When processing model files, the generic engine
171 parses parameters, variables, and mathematical expressions using a symbolic computing library, to
172 translate them into true functions. It also builds the simulation architecture and checks the consistency
173 of the model before running the simulation.

174 This separation between experts' and domain-specific knowledge (declarative part of the model)
175 on the one hand, and the algorithms to handle it (procedural part) on the other hand, is a classical, but
176 powerful Artificial Intelligence solution [45], known for helping experts to be involved directly in the
177 model design process, and for allowing fast, iterated feedbacks. Besides, this approach appears a kind
178 of "literate modelling" by analogy with Knuth's "literate programming" approach [46], aimed at
179 fostering a human-friendly, purpose-driven way of developing software codes. The elaboration of a DSL
180 for epidemiological models is actually a first attempt towards standardization, which must be supported
181 by an ability to encompass existing modelling paradigms and adapt to real use cases.

182 The modelling language defined in EMULSION is an "internal" DSL [47], as it is based on
183 another language, YAML (a human-friendly data serialization standard). Its syntax is quite simple,
184 relying mainly on lists and on dictionaries (i.e. key-value pairs), which can be nested one in another and,
185 for most components, do not require a special ordering. Contrary to most general-purpose programming
186 languages, this DSL is aimed at describing declarative knowledge, i.e. the model components and their
187 relations, the way to process them being implemented in the generic simulation engine. A whole
188 example is provided as supporting information (with syntactic colorization: Additional File, S4 Files).
189 Six main sections (first-level keys in the dictionary) must be specified: 1) the levels of interest in the

190 simulation (e.g. individuals, populations, metapopulations...) and their link to agent classes (i.e. either
191 agents defined in the generic simulation engine, as described below, or derived from the latter to provide
192 specific code add-ons); 2) the processes occurring at each level, which can be either handled through a
193 state machine (e.g. infection process, population dynamics, etc.), or implemented as a specific code add-
194 on in the class associated to the level; 3) the description of the state machines, composed of the list of
195 the states, with associated duration and actions if any, and the list of transitions between states, with
196 possible conditions and actions (the description of the state machines is equivalent to the state machine
197 diagram presented on Fig 3); 4) the comprehensive list of parameters used in the model, with their
198 description and values; 5) the list of agent variables (“statevars”) with their role; 6) the list of agent
199 actions with their description. Only the items of the two latter require subsequent implementation in the
200 proper agent classes as specific code add-ons. Additional features can be specified in the model, such
201 as time management (e.g. time unit, duration of discrete time steps, scheduling of events...) or desired
202 outputs.

203 The key point is that this description, which can be developed and consulted independently of
204 the generic simulation engine and of any possible code add-ons, does not require any computer science
205 skill to be understood and discussed. Hence, it fosters interactions with experts throughout modelling,
206 from formulating initial assumptions to specifying parameter values and identifying relevant scenarios
207 and outputs. Besides, an EMULSION model not only enforces an explicit specification of model items
208 that otherwise would be hidden in the code, but also requires a textual description of their rationale and
209 purpose, to keep track for instance of the evolution of assumptions or of the exact meaning of
210 parameters. Revising the model to account for new knowledge or to test alternative hypotheses
211 essentially consists in modifying the YAML file, by adding or removing states, transitions, parameters,
212 processes or actions, as shown on the application to Q Fever disease in the Results section.

213 **An agent-based software implementation**

214 Several elements in the model description file rely upon the agent-based software architecture
215 used in EMULSION, which is instantiated at runtime by the generic simulation to build the actual
216 simulation from the model description.

217 Multi-Agent Systems (MAS) [48, 49] have become a classical paradigm for the simulation of
218 complex systems. Agents, endowed with behaviors reflecting assumptions of a mechanistic model,

219 interact within a shared environment. They are quite flexible and can represent any kind of entity, since
 220 they are defined by their behaviors and interaction capabilities rather than by their structure. Their
 221 behaviors can be defined through rules, equations, probabilistic trials, etc. More recently, in multi-level
 222 MAS, agents are also used to explicitly represent intermediary organization levels within the system,
 223 with their own behavioral capabilities. Hence, they can be used to encapsulate other paradigms within
 224 a homogeneous interface. Among the few generic meta-models designed for multi-level agent-based
 225 simulation [50–53], we used the principles defined in [50] which proved flexible enough to adapt to
 226 other fields [54], and provide useful features for computational epidemiology, such as a clear separation
 227 between declarative and procedural aspects. A multi-level MAS, in this meta-model, is the combination
 228 of an architecture of nested agents (through a hosting relation) and an explicit and intelligible description
 229 of the processes where agents are involved.

230 State machines specify quite accurately individual behaviors for each possible process. In most
 231 situations however, keeping all individuals in the simulation would be highly inefficient and lack
 232 relevance. Therefore, agents can materialize groups at different levels, according to model requirements
 233 and to similarities between individuals. EMULSION intends to implement those principles, by
 234 combining agent classes defined to match typical relationships between a group and the underlying
 235 entities [55]. All agents are situated in at least one environment where they live, perceive, and act. Two
 236 main agent families are used: atoms representing individuals, and groups representing aggregation with
 237 a tunable granularity level.

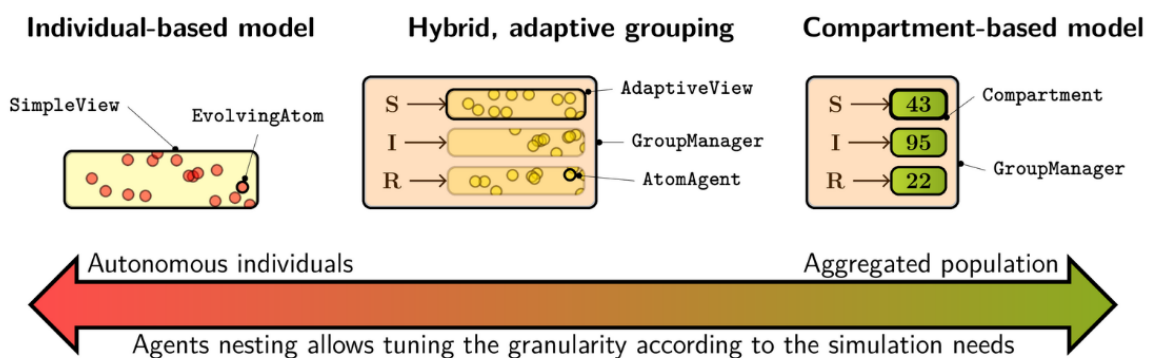


Fig 4. Integration of multiple modelling paradigms and scales. EMULSION allows multiple modelling paradigms to be expressed within the same formalism, based on nested agents, from the explicit representations of individuals (individual-based models) to aggregated populations (compartment-based models), with intermediary representations designed to group individuals depending on domain-dependent variables.

238 Depending on how such agents are combined and whether or not their behavior is controlled by
 239 a state machine, classical epidemiological modelling paradigms can be easily reproduced. A

240 “GroupManager” agent endowed with a health-related state machine and hosting several
241 “Compartment” agents reproduces a CBM (Fig 4, right). Contrarily, gathering “EvolvingAtom” agents,
242 each one owning a state machine, within a “SimpleView”, leads to an IBM (Fig 4, left). Refer to
243 Additional File, §S1 Appendix and S1 Fig for the detailed relationships between agent classes.

244 In addition, multi-level agents enable a hybrid modelling paradigm, mixing the preservation of
245 individual states as in IBM and the reduced computational cost of CBM. Indeed, aggregation structures
246 can be built at an intermediary stage, between individual- and population-oriented architectures.
247 Individuals can indeed be gathered according to each separate concern, based on their similarities
248 regarding each key variable. For instance, since the description of the infectious process is associated
249 with a specific state machine, each atom agent can be endowed with a variable for holding this state
250 (e.g. “health_state”). Then, it makes sense to gather individuals according to possible values of this
251 “health_state” variable, for instance using “AdaptiveView” agents, hosted by a “GroupManager” (Fig 4,
252 center). This “GroupManager” supervises the health-related state machine (instead of atoms) and, during
253 the simulation, determines how many atoms have to change their “health_state” value and move from
254 one “AdaptiveView” to another. To do so, due to the homogeneity of atoms within each
255 “AdaptiveView” regarding “health_state”, only one multinomial sample per group is required, instead
256 of one Bernoulli trial per individual, which reduces significantly the computation cost compared to a
257 classical IBM. Besides, using “AdaptiveView” agents as containers facilitates the separation of
258 concerns: if another process (e.g. recovery) suddenly affects the “health_state” variable of some agents,
259 their change is detected by the view, which asks its own host (the “GroupManager”) to move modified
260 atoms to the proper place.

261 Metapopulation appears a gathering of lower-level agents, such as those built after one of the
262 previous architectures. “MultiProcessManager” agents are designed to host them, provide a contact
263 structure, and be automatically constructed by EMULSION with the underlying components, based on
264 the specification of processes modelling the contact structure, key variables and state machines in the
265 model description file. Thus, several concerns are handled at the same time, without any special
266 development effort for the model designer.

267 EMULSION models can be checked prior to simulation to identify missing or inconsistent
268 information, and code templates can be generated to facilitate writing the specific add-ons. To run a

269 simulation from a model, EMULSION parses the YAML description file to read parameters, resolve
270 expressions, build the state machines, and instantiate the agent classes corresponding to the required
271 levels and groupings, based both on the objects of the generic simulation engine and on the specific code
272 add-ons.

273 **Application to the exploration of an epidemiological model (Q fever spread)**

274 The algorithms within EMULSION have been broadly tested on several very well-known
275 variants of SIR-like models [5, 6] based on CBM, IBM, hybrid modelling, and metapopulations.
276 However, the major added-value of EMULSION is to facilitate the development of complicated models
277 by model designers, to foster participative model revisions within short development time thanks to the
278 DSL. Hence, we addressed models for a real disease, Q fever in dairy cattle herds, for which herd
279 management processes have to be accounted for to reliably predict pathogen spread [15, 56, 57].

280 Q fever is a worldwide zoonosis caused by the bacterium *Coxiella burnetii*. It has recently
281 spread in Europe, e.g. in the Netherlands with a large number of human cases reported in 2007–2009
282 [58]. Domestic ruminants are recognized as the main source of human infection. In previous studies, a
283 detailed individual-based within-herd model was designed to help better controlling *C. burnetii* spread
284 in cattle herds with a particular attention paid to the diversity of transmission pathways and levels of
285 pathogen shedding by infected hosts [15]. The main parameters of a simplified variant of this model
286 were estimated from epidemiological data [56]. A study in the French department of Finistère revealed
287 seroprevalence levels for 2697 dairy herd by enzyme-linked immunosorbent assay (ELISA) in bulk tank
288 milk in 2012. 797 herds were detected seronegative in 2012 and tested again one year later. The annual
289 herd incidence (number of herds newly infected) was of 295 herds. The within-herd model was extended
290 to the between-herd level and confronted to such epidemiological data [57]. However, three main
291 computational and epidemiological issues remained. First, the infection process was mixed with the
292 reproduction cycle of cows, impeding modifications of biological assumptions and, thus, the exploration
293 of a larger variety of model structures. Second, the integration of within-herd dynamics into the between-
294 herd scale was not straightforward. Third, the simulated annual herd incidence was still too low.

295 We re-implemented these models using EMULSION to explore more quickly the interplay
296 between within-herd and between-herd levels. First, the original model [15] was simplified, keeping
297 relevant assumptions and removing those less crucial in the perspective of the between-herd dynamics

298 [59]. Then, it was extended to the regional between-herd level, where assumptions regarding airborne
299 transmission were compared. Finally, the within-herd model was revised with alternative hypotheses in
300 the infection process, to better reproduce the observed annual herd incidence at the between-herd level
301 while making plausible assumptions about host infection processes. The relative roles of trade and
302 airborne transmission in regional pathogen spread were reassessed under these new modelling
303 assumptions.

304 **Results**

305 **Model simplification within EMULSION**

306 According to the parsimony principle, we built a model with a minimum number of states,
307 transitions and parameters, trying to reproduce main simulation outcomes (prevalence, seroprevalence
308 and bacterial shedding) of the original within-herd model [15] after the transient period. To do so, we
309 identified possible simplifications (reduction of the number of states and transitions, and replacement
310 of distributions by aggregated parameters), and assessed them by simulation with a modified YAML
311 configuration file, without changing the code, and iterated the process with alternative hypotheses. In
312 the resulting model, called below "simplified model", 5 (out of 11) health states were retained
313 (Additional File, S2 Fig and S3 Fig): susceptible (S), infectious without (I⁻) or with (I⁺) antibodies, and
314 carrier with (C⁺) or without (C⁻) antibodies, without distinguishing between shedding levels or
315 pathways. Contaminations occurred through contacts with contaminated environment, bacteria (E_{total} ,
316 bacterial load in environment) coming either from local shedding (E_{local} , due to infectious animals from
317 the herd) or external sources (E_{aero} , from airborne transmission). At herd scale, when neglecting
318 between-herd transmission, E_{local} is equal to E_{total} (Additional File, S2 Appendix equation 1). The
319 probability of infection was determined by $p = 1 - e^{-\frac{N_0 * E_{total}}{N}}$, where N denoted the population
320 present in the herd (N_0 being a normalization factor). Besides, only adult female cows were taken into
321 account. They follow a reproduction process (Additional File, S4 Fig) with state depending on their
322 pregnancy status (P for pregnant, NP for non-pregnant). Transitions between P and NP were handled
323 by a duration in each state. Special events such as abortion could happen to pregnant cows up to three
324 weeks after infection. Local bacteria shedding occurred either during infectious states through on-stay

325 actions, or massively at special events (calving and abortion), through on-cross actions. The YAML file
326 corresponding to this model (model structure, levels, parameters, variables, etc.) is provided as
327 supporting information (Additional File, S4 Files).

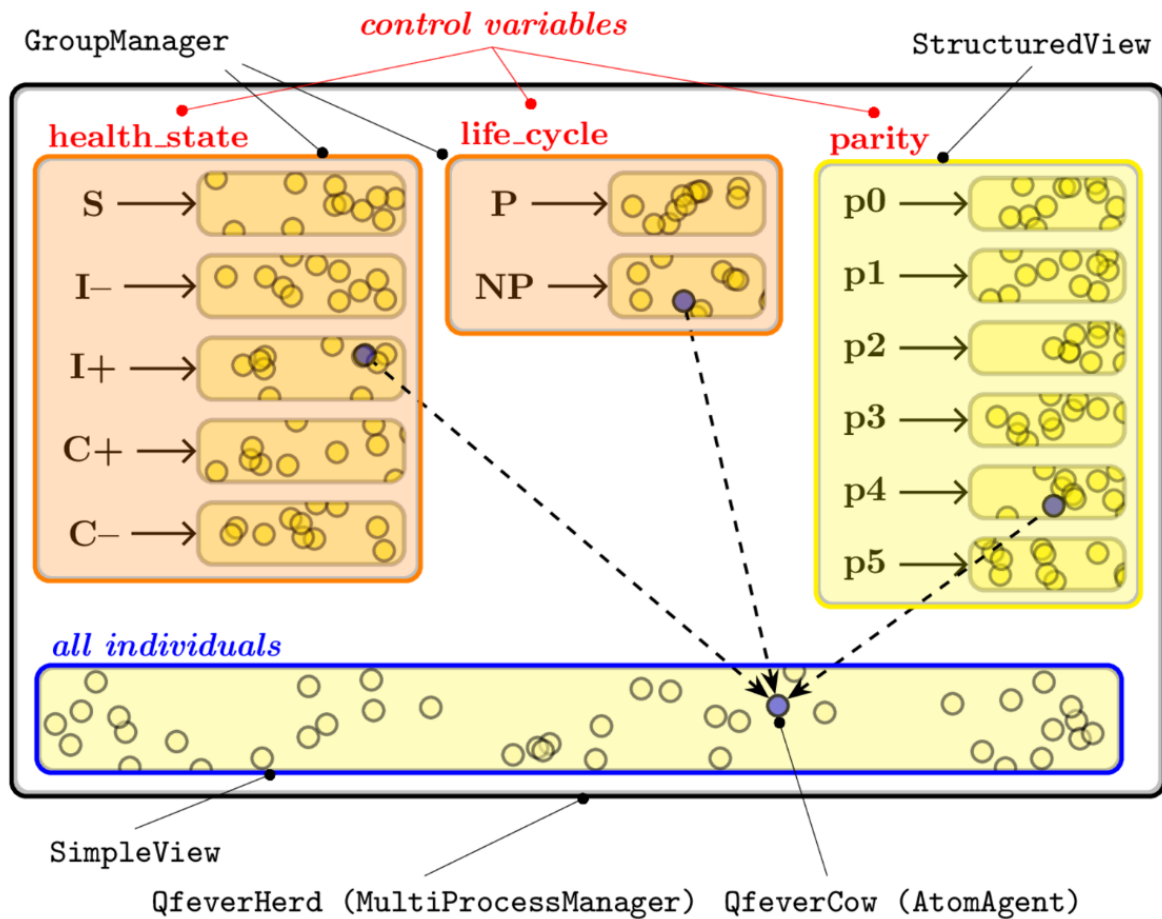


Fig 5. Structure of the within-herd hybrid model for Q fever. Individuals are aggregated according to concern-related variables (health state, life cycle, and parity). Individuals (e.g. the blue one in health state I+, life state NP, and parity 4) can be accessed through each concern or through a global list.

328 The model followed the hybrid structure (Fig 4, center) to fully account for individual events
329 (calving and abortion) without excessive computational cost. The implementation required only a class
330 for individuals, derived from "AtomAgent", and one for the herd, derived from "MultiProcessManager"
331 [59]. The resulting multi-agent architecture used to represent a herd is shown on Fig 5. Processes
332 involved at the herd level were the following: 1) culling (cow removal depending on parity); 2)
333 replacement (introduction of new animals); 3) infection and 4) farm management, both based on state
334 machines respectively affecting health state and life cycle; 5) actualization of animal grouping by parity;
335 6) bacterial decay in the environment (exponential decrease) and update of bacterial shedding. While
336 processes 1, 2, and 6 required writing short specific code add-ons, the others, involving generic
337 mechanisms such as state machines and groupings based on parity, health state, and life cycle, were

338 handled automatically by the generic simulation engine. Parameters were calibrated to match the median
339 and 10-90 percentiles (on 200 repetitions) of three major outcomes (prevalence, seroprevalence and
340 bacterial load in environment) of the original model simulation after the transient period (200 weeks
341 after introducing one I+ cow just before calving in a fully susceptible herd).

342 **Exploration: from within-herd to between-herd levels, back and forth**

343 Next, we focused on accounting for annual herd incidence observed in Finistère. The between-
344 herd model is a metapopulation, composed of independent herds linked through a contact network. As
345 in the within-herd architecture, the metapopulation could be implemented in EMULSION by a
346 “MultiProcessManager” agent, encapsulating a view holding all herds and endowed with dedicated
347 processes to handle interactions between them (Additional File, S7 Fig). Transforming the YAML file
348 from the within-herd to the between-herd scale only required to add the description architecture and
349 processes at the metapopulation level (Additional File, S1 Text), here assuming herds have similar
350 parameter values and sizes.

351 Herds could interact either through animal trade or airborne transmission from neighbor herds.
352 Initial assumptions considered animals bought from outside the metapopulation healthy. Bacteria were
353 transported and deposited by wind using a plume dispersion equation [60, 61] (called "Ermak-Stockie
354 function" below and detailed in Additional File, §S2 Appendix). Processes involved at the
355 metapopulation level thus were: 1) activation of herd processes; 2) airborne transmission; 3) selection
356 of animals for trade movements in source herds; 4) effective movement to destination herds. Herd
357 specificities (initial size, renewal, culling and trade movements) were based on the French livestock
358 exchange data, requiring specific code add-ons to make the metapopulation agent calibrate herd
359 parameters and extract the relevant trade movements from data. The predicted herd incidence with these
360 initial assumptions was much lower than in observed data (Fig 6-A). The discrepancy between observed
361 and predicted incidence could be explained either by a missing transmission route (but no other is known
362 for Q fever), wrong assumptions about risky trade or airborne transmission, or wrong assumptions about
363 the within-herd infection dynamics. To check for these two latter issues and improve herd incidence
364 predictions, we considered alternative assumptions on three main levers. We first assumed that animals
365 coming from outside the metapopulation had the same probability of being infected as inside the
366 metapopulation (rather than assuming them susceptible), this one being a part of a larger regional

367 population of herds. The impact on herd incidence was low (Fig 6-B) despite 18% of incoming
368 movements from outside the metapopulation. Second, regarding airborne transmission, we assumed that
369 individuals could be contaminated by inhalation of available bacteria (which is biologically plausible)
370 rather than by deposited pathogens. Hence, a simpler Gaussian approach [61], not accounting for
371 deposition, was implemented for plume dispersion (Additional File, equation 6, §S2 Appendix). Only a
372 few lines of codes were changed in the specific add-on to define the new function. Still, parameters
373 could not be calibrated within biologically plausible ranges to reach the expected seropositive herd
374 incidence (Fig 6-C), while the true herd incidence sharply increased. The difference between incidence
375 levels considering either infectious or seropositive animals, suggested to investigate the within-herd
376 model further. We reexamined shedding assumptions, shedding being observed to be intermittent [62].
377 The original model [15] assumed that I- cows were able to eliminate all bacteria and become S again
378 (non-shedder without antibodies and then apparently susceptible) [56], resulting in transitions from I-
379 to S. Alternatively, a latent state (L) could have been assumed, i.e. a non-shedding but infected state.
380 The intermittent shedding then can be explained by a loop between L and I- states (Additional File, S5
381 Fig), obviously increasing within-herd prevalence and reducing sharply spontaneous fade-out at local
382 scale. Using EMULSION, going back and forth from one model structure to another is straightforward,
383 even for a multi-scale model. The new within-herd model ("Latent state model") was built from the
384 Simplified model by adding a state and changing four transitions in the state machine describing health
385 states (i.e. 10 lines in the YAML file: Additional File, S2 Text). It was then calibrated to keep the same
386 steady-state regarding the same three main simulation outcomes (prevalence, seroprevalence and
387 bacterial load in environment) in the medium run as the simplified model (Additional File, S6 Fig).
388 Then, back to between-herd scale, we calibrated parameters in a plausible range for reaching the
389 expected herd incidence level (Fig 6-D).

390 A sensitivity analysis was carried out to assess the impact of model parameters on two main
391 outcomes at the metapopulation level (annual and weekly herd incidence), highlighting
392 N_0 (normalization factor), l (transition probability from L to I- states), m (transition probability from I-
393 to S), κ (proportion of bacteria leaving local environments to contribute to airborne transmission) and ζ
394 (contact rate with bacteria coming from airborne transmission) as key parameters (Additional File, S8

395 Fig). In addition, we ensured that spatial distributions of herds by health status were similar between
396 observed serological data and simulations (Additional File, S9 Fig).

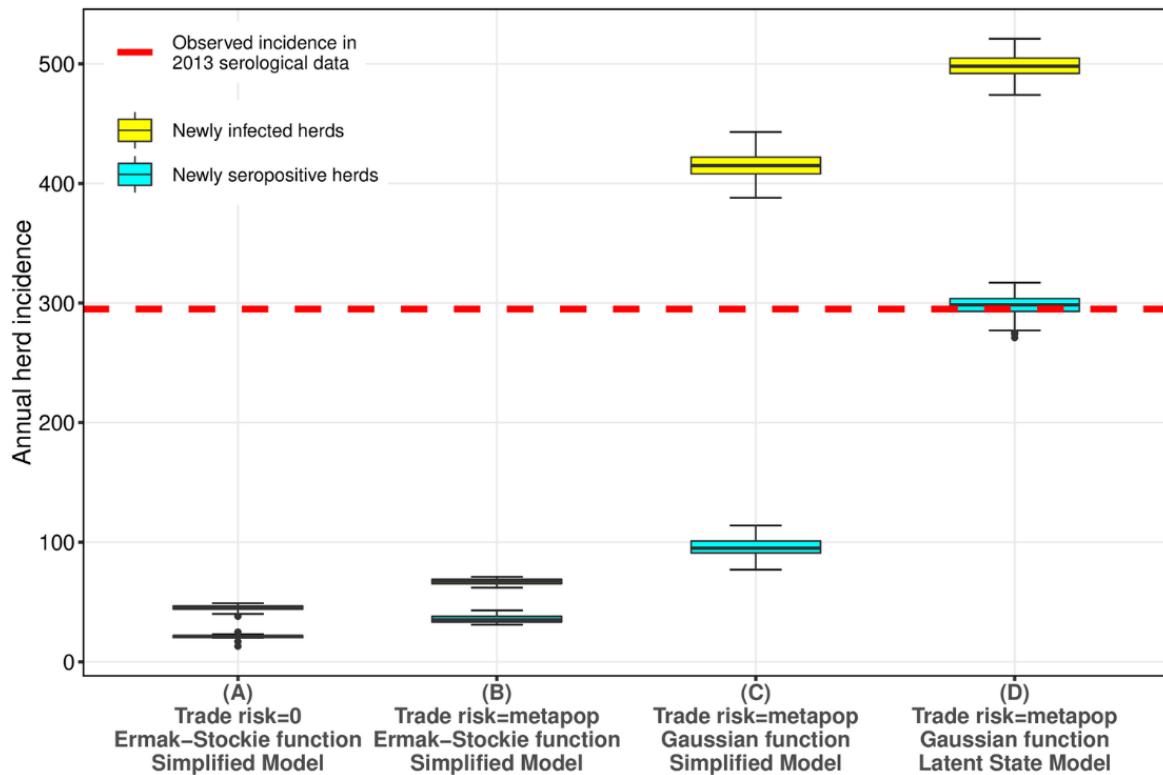


Fig 6. Annual herd incidence distributions in models based on combining variants of three main features. 1) No infection risk when purchasing animals from outside the metapopulation vs. risk similar to the average prevalence in the metapopulation, 2) Airborne transmission calculated by the Ermak-Stockie function vs. a Gaussian function, 3) Simplified within-herd model vs. model with a latent state. Yellow (cyan) color corresponds to the number of herds observed healthy (seronegative) in 2012 that hold a shedder (seropositive) animal at least once during the year. Observed data correspond to the 295 herds newly detected as seropositive by the ELISA test in 2013 (seronegative in 2012) in Finistère, France.

397 **Impact on previous conclusions**

398 Regarding pathways responsible for Q fever spread at the regional scale, we knew from the
399 previous study [57] that airborne transmission was predominant over trade movements. Yet, we wanted
400 to assess whether trade-borne infections could be neglected or not. After exploring model assumptions
401 and parameters to account for observed seroprevalence data in 2013 (Fig 6), we examined more finely
402 temporal and spatial effects of both transmission pathways, under the hypothesis that seroprevalence is
403 a relevant indicator for disease persistence within a one-year interval. First, infections of naive herds
404 appeared, as expected from previous results, to be caused mostly (89.5% of newly infected herds) by
405 airborne transmission (Fig 7, A). However, it also appeared that contaminations caused by trade
406 movements happened in a more deterministic way, which was not highlighted previously. When we
407 considered the infection dates of herds contaminated by trade movements at least once in 50 stochastic

408 repetitions by trade movements, two groups were identified (Fig 7, B): first, herds infected early in the
409 simulations, almost always at the same date and by trade movements; second, herds infected at a variable
410 date and with a variable contribution of airborne transmission. Spatial distribution of incident herds
411 (Fig 7, C) pointed out that areas with a high density of initially infected herds (and of herds in general:
412 Additional File, S10 Fig) drove at the same time the predominance of airborne transmission and the
413 probability that a herd subject to airborne transmission risk becomes infected. Conversely, herds infected
414 at least once by trade movements were mostly located on peripheral areas, such as coasts (Fig 7, D), and
415 most of them purchased animals directly in an initially prevalent herd before becoming infected, with
416 mostly early infection dates. To summarize, the spread of Q fever within areas of high prevalence and
417 high herd density is strong and mainly caused by airborne transmission, which argues for vaccination
418 as a disease control strategy in such areas, while herds in low prevalence areas have little chance to be
419 contaminated but by trade, which supports tests on purchase in that case.

420 **Discussion and conclusions**

421 EMULSION is the first framework that simultaneously provides a Domain-Specific Language
422 dedicated to the comprehensive and accurate description of epidemiological models, from SIR-like
423 models to more complex multi-scale multi-concern models, together with a modular simulation engine
424 using a multi-level agent-based architecture, to encompass existing epidemiological issues and
425 modelling paradigms within a homogeneous interface (Fig 1).

426 Elaborating realistic models (such as [15, 57]) often requires many trials aiming at the
427 exploration of various assumptions and processes. EMULSION significantly accelerates model
428 development, first because it provides classical computational bricks, but also since changing
429 hypotheses (e.g. adding or deleting a state or a transition) generally consists in modifying the
430 configuration file, instead of rewriting many parts of a large specific source code. The modularity of
431 model description allows assessing separately hypotheses, which demands much more work when
432 dealing with ad-hoc models and is more prone to programming and hence interpretation errors.

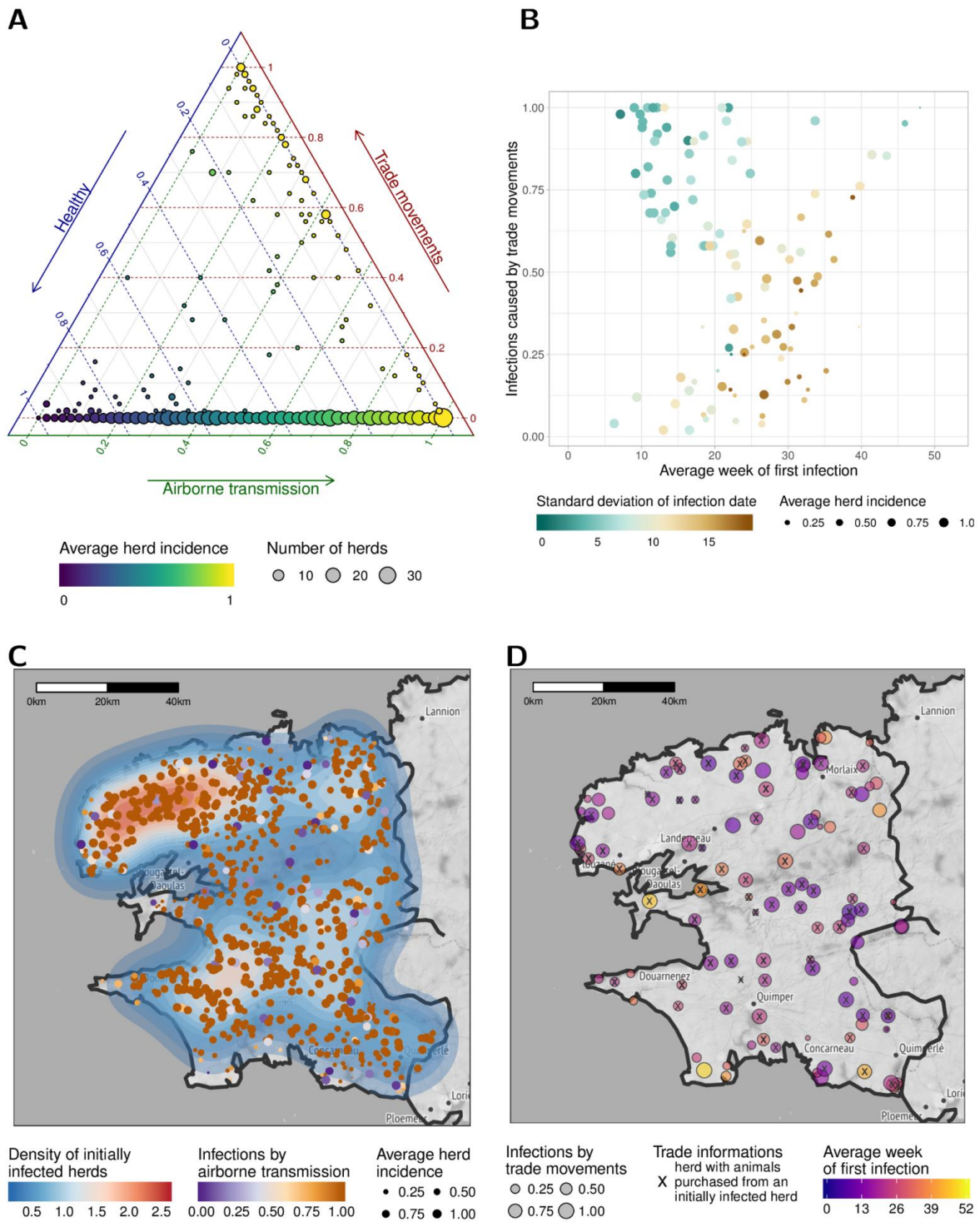


Fig 7. Contributions of infection pathways at regional scale, as predicted by the between-herd Q fever model with latent state and a Gaussian plume airborne transmission (50 stochastic repetitions of the standard scenario). (A) Distribution of the proportion of repetitions where each herd stayed healthy, became infected by animal trade movements, or became infected by airborne transmission, over one year (color shows the proportion of repetitions where the herd became infected by any of the transmission routes). (B) Amongst herds infected in at least one repetition by trade movements, relation between the proportion of infections caused by movement and the infection date (in average and standard deviation), exhibiting two subgroups: one with an early and little variable infection date, caused very often by movements (blue points), and the other with a more variable infection date and caused less often by movements (brown points). (C) Map of Finistère with the density of initially infected herds (2012 serological data) and the location of herds infected in at least one repetition. Color shows the proportion of infections caused by airborne transmission vs. trade movements. (D) Map showing the location and average infection date of herds infected at least once by trade movements. Herds marked with a "x" purchased animals from initially infected herds before their own infection. Map tiles by Stamen Design, under CC BY 3.0. Data by OpenStreetMap, under ODbL.

433 Using EMULSION, we re-implemented very quickly the initial within-herd Q fever dairy cattle
434 model and assessed the compatibility of simpler alternative model structures with previous predictions
435 derived from a model whose parameters were estimated using observed on-farm data. Then, moving to
436 the regional scale required little transformation of the within-herd model. Small pieces of code related
437 to trade and weather data and calculation of transmission functions were used as specific add-ons by the
438 simulation engine. The facility to compare a large panel of assumptions regarding the existence of a
439 latent stage in Q fever infectious process, the nature of airborne transmission, and values of the less
440 well-known parameters, allowed us to identify and calibrate the best candidates with respect to observed
441 herd incidence data. Also, conditions under which trade movements and airborne transmission
442 contribute to new infections were explored, highlighting new findings, especially that infections caused
443 by movements are almost deterministic and impact mostly herds in peripheral areas and low prevalence
444 areas.

445 This work provides a proof of concept, demonstrating the added-value of using such a
446 framework, both in terms of code reduction and model readability. To promote our approach and
447 modelling language, EMULSION will be soon released as an open-source software. The current version
448 of the generic simulation engine being written in Python, efficiency cannot compete with compiled
449 languages such as C++ (as the code generated by KENDRICK [40]) or Java (used in the Broadwick
450 framework [38]), especially at the between-herd scale. This was not under the scope of the present study.
451 Nevertheless, the facility to choose the granularity level of simulations and the adaptive gathering of
452 individuals make the approach much more efficient than IBM anyway. To tackle efficiency issue, the
453 next step will be to consider using EMULSION's DSL to build dedicated, optimized code from model
454 descriptions.

455 We showed how modelling paradigms and scales could be wrapped within agents, which are in
456 charge of processing the required calculations according to their specificities. This allows modellers to
457 focus on their research questions instead of implementation issues, while still being able to select and
458 compare relevant modeling paradigms. To go further, new computational issues in epidemiological
459 modelling have to be addressed, especially in coupling contrasted paradigms [63]. For instance, multi-
460 host pathosystems may combine populations having highly contrasted characteristics, such as size (a
461 large size leading to more deterministic dynamics while a small one enhances stochastic events) and

462 movement patterns (vectors spread continuously while livestock trade and human activity give rise to
463 discrete long-distance jumps). In addition, the level of required details may differ, possibly leading to
464 combine aggregated representation (i.e. CBM) with preservation of individuals (i.e. IBM). The
465 integration of such features into both the DSL and the generic engine (especially as new agent classes)
466 will enable modellers to address such computational challenges in epidemiological modelling.

467 From the point of view of computer science, adapting the original multi-level agent-based meta-
468 model [41, 50] to epidemiological issues also was fruitful. Dealing with aggregation and disaggregation
469 in an adaptive way is a challenging open question in multi-level modelling. The architecture designed
470 for multi-scale epidemiological systems will provide clues for building similar structures in general-
471 purpose agent-based systems, and continue the identification and characterization of design patterns in
472 multi-level agent-based simulation initiated in [55].

473 We consider our contribution a first step towards a standardized DSL for epidemiology. Though
474 initiated in the context of animal health, our approach is generic and modular enough to extend to human
475 and plant epidemiology. The generalization of such methods could enhance significantly the
476 reactivity of modelers in sketching, assessing, and recommending reliable and efficient control
477 measures against outbreaks, accounting for possible biases in model predictions arising from uncertainty
478 in model assumptions.

479 **References**

- 480 1. Peng RD. Reproducible Epidemiologic Research. *Am J Epidemiol.* 2006;163:783–789.
- 481 2. Sandve GK, Nekrutenko A, Taylor J, Hovig E. Ten Simple Rules for Reproducible Computational Research.
482 *PLoS Comput Biol.* 2013;9:e1003285.
- 483 3. Leek JT, Peng RD. Opinion: Reproducible research can still be wrong: Adopting a prevention approach. *Proc*
484 *Natl Acad Sci.* 2015;112:1645–1646.
- 485 4. Kermack WO, McKendrick AG. A Contribution to the Mathematical Theory of Epidemics. *Proc R Soc.*
486 *1927;;700–721.*
- 487 5. Diekmann O, Heesterbeek HJ. *Mathematical epidemiology of infectious diseases: model buildign, analysis and*
488 *interpretation.* Chichester: Wiley; 2000.
- 489 6. Keeling MJ, Rohani P. *Modeling Infectious Diseases in Humans and Animals.* Princeton University Press; 2008.
- 490 7. Merali Z. Computational science: ...Error. *Nature.* 2010;467:775–777.
- 491 8. Hethcote HW. *The Mathematics of Infectious Diseases.* SIAM Rev. 2000;42:599–653.

- 492 9. Marcé C, Ezanno P, Seegers H, Pfeiffer D, Fourichon C. Predicting fadeout versus persistence of
493 paratuberculosis in a dairy cattle herd for management and control purposes: a modelling study. *Vet Res.*
494 2011;42:36.
- 495 10. DeAngelis DL, Grimm V. Individual-based models in ecology after four decades. *F1000Prime Rep.* 2014;6.
496 doi:10.12703/P6-39.
- 497 11. Railsback SF, Grimm V. *Agent-Based and Individual-Based Modelling: A Practical Introduction.* Princeton
498 University Press; 2011.
- 499 12. Ferguson NM, Cummings DAT, Cauchemez S, Fraser C, Riley S, Meeyai A, et al. Strategies for containing
500 an emerging influenza pandemic in Southeast Asia. *Nature.* 2005;437:209–14.
- 501 13. Halloran ME, Ferguson NM, Eubank S, Longini IM, Cummings DAT, Lewis B, et al. Modeling targeted
502 layered containment of an influenza pandemic in the United States. *Proc Natl Acad Sci.* 2008;105:4639–44.
- 503 14. Amouroux E, Desvaux S, Drogoul A. Towards Virtual Epidemiology: An Agent-Based Approach to the
504 Modeling of H5N1 Propagation and Persistence in North-Vietnam. In: Bui TD, Ho TV, Ha QT, editors. 11th
505 Pacific Rim Int. Conf. on Multi-Agents (PRIMA). Springer; 2008. p. 26–33. doi:10.1007/978-3-540-89674-
506 6_6.
- 507 15. Courcoul A, Monod H, Nielen M, Klinkenberg D, Hogerwerf L, Beaudeau F, et al. Modelling the effect of
508 heterogeneity of shedding on the within herd *Coxiella burnetii* spread and identification of key parameters
509 by sensitivity analysis. *J Theor Biol.* 2011;284:130–141.
- 510 16. Robins J, Bogen S, Francis A, Westhoek A, Kanarek A, Lenhart S, et al. Agent-based model for Johne’s disease
511 dynamics in a dairy herd. *Vet Res.* 2015;46. doi:10.1186/s13567-015-0195-y.
- 512 17. Marshall BDL, Galea S. Formalizing the Role of Agent-Based Modeling in Causal Inference and
513 Epidemiology. *Am J Epidemiol.* 2014;181:92–99.
- 514 18. Parker J, Epstein JM. A Distributed Platform for Global-Scale Agent-Based Models of Disease Transmission.
515 *ACM Trans Model Comput Simul.* 2011;22:2:1–2:25.
- 516 19. Gilpin M, Hanski I, editors. *Metapopulation Dynamics: Empirical and Theoretical Investigations.* Elsevier BV;
517 1991. doi:10.1016/b978-0-12-284120-0.50003-6.
- 518 20. Grenfell B, Harwood J. (Meta)population dynamics of infectious diseases. *Trends Ecol Evol.* 1997;12:395–
519 399.
- 520 21. Keeling M. The implications of network structure for epidemic dynamics. *Theor Popul Biol.* 2005;67:1–8.
- 521 22. Arino J, van den Driessche P. Disease spread in metapopulations. In: Brunner H, Zhao X-Q, Zou X, editors.
522 *Nonlinear Dynamics and Evolution Equations.* American Mathematical Society; 2006. p. 1–12.
- 523 23. Beaunée G, Vergu E, Ezanno P. Modelling of paratuberculosis spread between dairy cattle farms at a regional
524 scale. *Vet Res.* 2015;46. doi:10.1186/s13567-015-0247-3.
- 525 24. Ajelli M, Gonçalves B, Balcan D, Colizza V, Hu H, Ramasco JJ, et al. Comparing large-scale computational
526 approaches to epidemic modeling: Agent-based versus structured metapopulation models. *BMC Infect Dis.*
527 2010;10. doi:10.1186/1471-2334-10-190.
- 528 25. Keeling MJ, Danon L, Vernon MC, House TA. Individual identity and movement networks for disease
529 metapopulations. *Proc Natl Acad Sci.* 2010;107:8866–8870.
- 530 26. Gillespie DT. Exact stochastic simulation of coupled chemical reactions. *J Phys Chem.* 1977;81:2340–61.
- 531 27. Bretó C, He D, Ionides EL, King AA. Time series analysis via mechanistic models. *Ann Appl Stat.*
532 2009;3:319–48.

- 533 28. Le Novere N, Hucka M, Mi H, Moodie S, Schreiber F, Sorokin A, et al. The Systems Biology Graphical
534 Notation. *Nat Biotech.* 2009;27:735–741.
- 535 29. Díaz-Zuccarini V, Pichardo-Almarza C. On the formalization of multi-scale and multi-science processes for
536 integrative biology. *Interface Focus.* 2011;1:426–437.
- 537 30. Perra N, Balcan D, Gonçalves B, Vespignani A. Towards a Characterization of Behavior-Disease Models.
538 *PLoS ONE.* 2011;6:e23084.
- 539 31. Grimm V, Berger U, Bastiansen F, Eliassen S, Ginot V, Giske J, et al. A standard protocol for describing
540 individual-based and agent-based models. *Ecol Model.* 2006;198:115–126.
- 541 32. Amouroux E, Gaudou B, Desvaux S, Drogoul A. O.D.D.: A Promising but Incomplete Formalism for
542 Individual-Based Model Specification. In: *RIVF Int. Conf. on Computing and Communication*
543 *Technologies.* IEEE; 2010. doi:10.1109/rivf.2010.5633421.
- 544 33. Roche B, Guégan J-F, Bousquet F. Multi-agent systems in epidemiology: a first step for computational biology
545 in the study of vector-borne disease transmission. *BMC Bioinformatics.* 2008;9. doi:10.1186/1471-2105-9-
546 435.
- 547 34. Collier N, Ozik J, Macal CM. Large-Scale Agent-Based Modeling with Repast HPC: A Case Study in
548 Parallelizing an Agent-Based Model. In: *Parallel Processing Workshops (Euro-Par).* Springer Nature; 2015.
549 p. 454–465. doi:10.1007/978-3-319-27308-2_37.
- 550 35. Widgren S, Bauer P, Engblom S. SimInf: An R package for Data-driven Stochastic Disease Spread
551 Simulations. *ArXiv Prepr ArXiv160501421 Q-BioPE.* 2016. <http://arxiv.org/abs/1605.01421>.
- 552 36. Cakici B, Boman M. A workflow for software development within computational epidemiology. *J Comput*
553 *Sci.* 2011;2:216–222.
- 554 37. Broeck WV den, Gioannini C, Gonçalves B, Quaggiotto M, Colizza V, Vespignani A. The GLEaMviz
555 computational tool, a publicly available software to explore realistic epidemic spreading scenarios at the
556 global scale. *BMC Infect Dis.* 2011;11. doi:10.1186/1471-2334-11-37.
- 557 38. O’Hare A, Lycett SJ, Doherty T, Salvador LCM, Kao RR. Broadwick: a framework for computational
558 epidemiology. *BMC Bioinformatics.* 2016;17. doi:10.1186/s12859-016-0903-2.
- 559 39. Haddad H, Moulin B, Thériault M. A fully GIS-integrated simulation approach for analyzing the spread of
560 epidemics in urban areas. *SIGSPATIAL Spec.* 2016;8:34–41.
- 561 40. Bui T-M-A, Ziane M, Stinckwich S, Ho T-V, Roche B, Papoulias N. Separation of Concerns in
562 Epidemiological Modelling. In: Fuentes L, Batory DS, Czarnecki K, editors. *Proceedings of the 15th*
563 *International Conference on Modularity.* ACM; 2016. p. 196–200. doi:10.1145/2892664.2892699.
- 564 41. Picault S, Huang Y-L, Sicard V, Ezanno P. Enhancing Sustainability of Complex Epidemiological Models
565 through a Generic Multilevel Agent-based Approach. In: Sierra C, editor. *Proceedings of the 26th*
566 *International Joint Conference on Artificial Intelligence (IJCAI’2017).* Melbourne, Australia: AAAI; 2017.
- 567 42. Bobashev GV, Goedecke DM, Yu F, Epstein JM. A Hybrid Epidemic Model: Combining The Advantages Of
568 Agent-Based And Equation-Based Approaches. *Winter Simul Conf.* 2007. doi:10.1109/wsc.2007.4419767.
- 569 43. Booth TL. *Sequential Machines and Automata Theory.* 1st edition. New York: John Wiley and Sons; 1967.
- 570 44. Mernik M, Heering J, Sloane AM. When and how to develop domain-specific languages. *ACM Comput Surv.*
571 2005;37:316–44.
- 572 45. Newell A, Shaw JC, Simon HA. Report on a General Problem-Solver Program. In: *Proceedings of the*
573 *International Conference on Information Processing.* 1959. p. 256–264.
- 574 46. Knuth DE. *Literate programming.* Stanford, Calif.: Center for the Study of Language and Information; 1992.

- 575 47. Fowler M, Parsons R. Domain-specific languages. Upper Saddle River, NJ: Addison-Wesley; 2011.
- 576 48. Ferber J. Multi-agent systems: an introduction to distributed artificial intelligence. Harlow: Addison-Wesley;
- 577 1998.
- 578 49. Weiss G, editor. Multiagent systems: a modern approach to distributed artificial intelligence. Cambridge, Mass:
- 579 MIT Press; 1999.
- 580 50. Picault S, Mathieu P. An Interaction-Oriented Model for Multi-Scale Simulation. In: Walsh T, editor.
- 581 Proceedings of the 22nd International Joint Conference on Artificial Intelligence (IJCAI'2011). AAAI; 2011.
- 582 p. 332–337. <https://hal.archives-ouvertes.fr/hal-00826401>.
- 583 51. Morvan G, Veremme A, Dupont D. IRM4MLS: The Influence Reaction Model for Multi-Level Simulation.
- 584 In: Multi-Agent-Based Simulation XI. Springer; 2011. p. 16–27. doi:10.1007/978-3-642-18345-4_2.
- 585 52. Camus B, Bourjot C, Chevrier V. Multi-level modeling as a society of interacting models. In: Yilmaz L, Ören
- 586 TI, Madey G, Sierhuis M, Zhang Y, editors. Agent-Directed Simulation Symposium (in SpringSim).
- 587 SCS/ACM; 2013. <http://dl.acm.org/citation.cfm?id=2499595>.
- 588 53. Huraux T, Sabouret N, Haradji Y. A Multi-level Model for Multi-agent based Simulation: In: Proceedings of
- 589 the 6th International Conference on Agents and Artificial Intelligence. SCITEPRESS - Science and and
- 590 Technology Publications; 2014. p. 139–46. doi:10.5220/0004814501390146.
- 591 54. Maudet A, Touya G, Duchêne C, Picault S. DIOGEN, a multi-level oriented model for cartographic
- 592 generalization. Int J Cartogr. 2017;3:121–33.
- 593 55. Mathieu P, Morvan G, Picault S. Multi-level agent-based simulations: Four design patterns. Simul Model Pract
- 594 Theory. 2018;in press.
- 595 56. Courcoul A, Vergu E, Denis J-B, Beaudeau F. Spread of Q fever within dairy cattle herds: key parameters
- 596 inferred using a Bayesian approach. Proc R Soc B Biol Sci. 2010;277:2857–65.
- 597 57. Pandit P, Hoch T, Ezanno P, Beaudeau F, Vergu E. Spread of *Coxiella burnetii* between dairy cattle herds in
- 598 an enzootic region: modelling contributions of airborne transmission and trade. Vet Res. 2016;47.
- 599 doi:10.1186/s13567-016-0330-4.
- 600 58. van der Hoek W, Morroy G, Renders NHM, Wever PC, Hermans MHA, Leenders ACAP, et al. Epidemic Q
- 601 Fever in Humans in the Netherlands. In: Toman R, Heinzen RA, Samuel JE, Mege J-L, editors. *Coxiella*
- 602 *burnetii*: Recent Advances and New Perspectives in Research of the Q Fever Bacterium. Dordrecht: Springer
- 603 Netherlands; 2012. p. 329–64. doi:10.1007/978-94-007-4315-1_17.
- 604 59. Picault S, Huang Y-L, Sicard V, Beaudeau F, Ezanno P. A Multi-Level Multi-Agent Simulation Framework
- 605 in Animal Epidemiology. In: Demazeau Y, Davidsson P, Vale Z, Bajo J, editors. Proceedings of the 15th
- 606 International Conference on Practical Applications of Agents and Multi-Agent Systems (PAAMS'2017).
- 607 Porto: Springer; 2017. p. 209–21.
- 608 60. Ermak DL. An analytical model for air pollutant transport and deposition from a point source. Atmospheric
- 609 Environ 1967. 1977;11:231–7.
- 610 61. Stockie JM. The Mathematics of Atmospheric Dispersion Modeling. SIAM Rev. 2011;53:349–72.
- 611 62. Guatteo R, Beaudeau F, Joly A, Seegers H. *Coxiella burnetii* shedding by dairy cows. Vet Res. 2007;38:849–
- 612 60.
- 613 63. Sutherland WJ, Freckleton RP, Godfray H CJ, Beissinger SR, Benton T, Cameron DD, et al. Identification of
- 614 100 fundamental ecological questions. J Ecol. 2013;101:58–67.
- 615

616 **Additional material**

617 **Additional File. (PDF) Complementary information for the main article.** Contains a description of agent
618 classes, of airborne transmission functions, additional figures, and YAML file descriptions or modifications.

619 **Abbreviations**

620 **CBM:** Compartment-Based Model; **DSL:** Domain-Specific Language; **IBM:** Individual-Based Model; **MAS:**
621 Multi-Agent Systems; **ODD:** "Overview, Design concepts, Details" protocol; **ODE:** Ordinary Differential
622 Equations; **SBGN:** Systems Biology Graphical Notation

623 **Declarations**

624 **Ethics approval and consent to participate.** Not applicable.

625 **Consent for publication.** Not applicable.

626 **Competing interests.** The authors declare they have no competing interests.

627 **Availability of data and material. Data:** French livestock exchange data were provided by the French Ministry
628 of Agriculture (FMA). Data collection and analyses are subject to a confidentiality agreement (available upon
629 request from the following contact point: bicma.sdspa.dgal@agriculture.gouv.fr). Public weather datasets were
630 provided by the European Centre for Medium-Range Weather Forecasts (ECMWF). Parameters and structure of
631 Q Fever models are fully available in Supporting Information. **Software:** The framework EMULSION is awaiting
632 approval from the French National Institute for Agricultural Research (INRA) before being publicly released. In
633 the meanwhile, the specific code add-ons used for Q fever modelling and the generic simulation engine are
634 available upon request from the contact author.

635 **Funding.** The work was funded by the French Research Agency (ANR) through projects MIHMES (ANR-10-
636 BINF-07) and CADENCE (ANR-16-CE32-0007-01), the European fund for the Regional Development (FEDER)
637 of Pays-de-la-Loire, and the Animal Health Division of INRA.

638 **Authors' contributions.** SP designed the DSL and the engine of EMULSION, led software developments, and
639 drafted the manuscript; YLH participated in software development, carried out Q fever study, and drafted the
640 corresponding section; VS participated in the development and data preparation; TH, EV and FB provided
641 expertise on Q fever and participated in results interpretation; PE supervised the specifications of EMULSION,
642 and designed the Q fever study. All authors helped draft the manuscript and gave final approval for publication.

643 **Acknowledgments.** We are grateful to the French Ministry of Agriculture for granting us access to cattle datasets.